

# Calculus, Applications and Theory

Kenneth Kuttler

April 29, 2004



# Contents

<b>1</b>	<b>Introduction</b>	<b>11</b>
<b>I</b>	<b>Preliminaries</b>	<b>13</b>
<b>2</b>	<b>The Real Numbers</b>	<b>15</b>
2.1	The Number Line And Algebra Of The Real Numbers . . . . .	15
2.2	Exercises . . . . .	19
2.3	Order . . . . .	20
2.3.1	Set Notation . . . . .	21
2.4	Exercises With Answers . . . . .	22
2.5	Exercises . . . . .	23
2.6	The Absolute Value . . . . .	23
2.7	Exercises . . . . .	26
2.8	Well Ordering Principle And Archimedian Property . . . . .	27
2.9	Exercises . . . . .	30
2.10	Divisibility And The Fundamental Theorem Of Arithmetic . . . . .	33
2.11	Exercises . . . . .	35
2.12	Systems Of Equations . . . . .	36
2.13	Exercises . . . . .	40
2.14	Completeness of $\mathbb{R}$ . . . . .	42
2.15	Review Exercises . . . . .	43
<b>3</b>	<b>Basic Geometry And Trigonometry</b>	<b>47</b>
3.1	Similar Triangles And Pythagorean Theorem . . . . .	47
3.2	Cartesian Coordinates And Straight Lines . . . . .	49
3.3	Exercises . . . . .	52
3.4	Distance Formula And Trigonometric Functions . . . . .	52
3.5	The Circular Arc Subtended By An Angle . . . . .	55
3.6	The Trigonometric Functions . . . . .	60
3.7	Exercises . . . . .	63
3.8	Some Basic Area Formulas . . . . .	65
3.8.1	Areas Of Triangles And Parallelograms . . . . .	65
3.8.2	The Area Of A Circular Sector . . . . .	66
3.9	Exercises . . . . .	67
3.10	Parabolas, Ellipses, and Hyperbolas . . . . .	69
3.10.1	The Parabola . . . . .	69
3.10.2	The Ellipse . . . . .	70
3.10.3	The Hyperbola . . . . .	71

3.11 Exercises . . . . .	72
<b>4 The Complex Numbers</b>	<b>73</b>
4.1 Exercises . . . . .	76
 <b>II Functions Of One Variable</b>	 <b>79</b>
<b>5 Functions</b>	<b>81</b>
5.1 General Considerations . . . . .	81
5.2 Exercises . . . . .	85
5.3 Continuous Functions . . . . .	87
5.4 Sufficient Conditions For Continuity . . . . .	91
5.5 Continuity Of Circular Functions . . . . .	92
5.6 Exercises . . . . .	93
5.7 Properties Of Continuous Functions . . . . .	93
5.8 Exercises . . . . .	97
5.9 Limits Of A Function . . . . .	97
5.10 Exercises . . . . .	102
5.11 The Limit Of A Sequence . . . . .	103
5.11.1 Sequences And Completeness . . . . .	106
5.11.2 Decimals . . . . .	109
5.11.3 Continuity And The Limit Of A Sequence . . . . .	110
5.12 Exercises . . . . .	110
5.13 Uniform Continuity . . . . .	113
5.14 Exercises . . . . .	114
5.15 Theorems About Continuous Functions . . . . .	114
 <b>6 Derivatives</b>	 <b>119</b>
6.1 Velocity . . . . .	119
6.2 The Derivative . . . . .	120
6.3 Exercises With Answers . . . . .	125
6.4 Exercises . . . . .	126
6.5 Local Extrema . . . . .	126
6.6 Exercises With Answers . . . . .	128
6.7 Exercises . . . . .	130
6.8 Mean Value Theorem . . . . .	132
6.9 Exercises . . . . .	134
6.10 Curve Sketching . . . . .	135
6.11 Exercises . . . . .	136
 <b>7 Some Important Special Functions</b>	 <b>139</b>
7.1 The Circular Functions . . . . .	139
7.2 Exercises . . . . .	141
7.3 The Exponential And Log Functions . . . . .	142
7.3.1 The Rules Of Exponents . . . . .	142
7.3.2 The Exponential Functions, A Wild Assumption . . . . .	143
7.3.3 The Special Number, $e$ . . . . .	145
7.3.4 The Function $\ln  x $ . . . . .	146
7.3.5 Logarithm Functions . . . . .	146
7.4 Exercises . . . . .	147

<b>8</b>	<b>Properties And Applications Of Derivatives</b>	<b>149</b>
8.1	The Chain Rule And Derivatives Of Inverse Functions . . . . .	149
8.1.1	The Chain Rule . . . . .	149
8.1.2	Implicit Differentiation And Derivatives Of Inverse Functions . . . . .	150
8.2	Exercises . . . . .	152
8.3	The Function $x^r$ For $r$ A Real Number . . . . .	153
8.3.1	Logarithmic Differentiation . . . . .	154
8.4	Exercises . . . . .	155
8.5	The Inverse Trigonometric Functions . . . . .	156
8.6	The Hyperbolic And Inverse Hyperbolic Functions . . . . .	159
8.7	Exercises . . . . .	160
8.8	L'Hôpital's Rule . . . . .	161
8.8.1	Interest Compounded Continuously . . . . .	165
8.9	Exercises . . . . .	166
8.10	Related Rates . . . . .	167
8.11	Exercises . . . . .	168
8.12	The Derivative And Optimization . . . . .	170
8.13	Exercises . . . . .	173
8.14	The Newton Raphson Method . . . . .	175
8.15	Exercises . . . . .	177
<b>9</b>	<b>Antiderivatives And Differential Equations</b>	<b>179</b>
9.1	Initial Value Problems . . . . .	179
9.2	Areas . . . . .	181
9.3	Area Between Graphs . . . . .	182
9.4	Exercises . . . . .	184
9.5	The Method Of Substitution . . . . .	186
9.6	Exercises . . . . .	188
9.7	Integration By Parts . . . . .	190
9.8	Exercises . . . . .	191
9.9	Trig. Substitutions . . . . .	192
9.10	Exercises . . . . .	196
9.11	Partial Fractions . . . . .	197
9.12	Rational Functions Of Trig. Functions . . . . .	202
9.13	Exercises . . . . .	202
9.14	Practice Problems For Antiderivatives . . . . .	203
9.15	Volumes . . . . .	210
9.16	Exercises . . . . .	213
9.17	Lengths And Areas Of Surfaces Of Revolution . . . . .	214
9.18	Exercises . . . . .	219
9.19	The Equation $y' + ay = 0$ . . . . .	220
9.20	Exercises . . . . .	221
9.21	Force On A Dam And Work Of A Pump . . . . .	222
9.22	Exercises . . . . .	224
<b>10</b>	<b>The Integral</b>	<b>227</b>
10.1	Upper And Lower Sums . . . . .	227
10.2	Exercises . . . . .	230
10.3	Functions Of Riemann Integrable Functions . . . . .	231
10.4	Properties Of The Integral . . . . .	233
10.5	Fundamental Theorem Of Calculus . . . . .	237

10.6 Exercises . . . . .	240
10.7 Return Of The Wild Assumption . . . . .	241
10.8 Exercises . . . . .	245
10.9 Techniques Of Integration . . . . .	247
10.9.1 The Method Of Substitution . . . . .	247
10.9.2 Integration By Parts . . . . .	248
10.10 Exercises . . . . .	249
10.11 Improper Integrals . . . . .	252
10.12 Exercises . . . . .	255
<b>11 Infinite Series</b>	<b>257</b>
11.1 Approximation By Taylor Polynomials . . . . .	257
11.2 Exercises . . . . .	259
11.3 Infinite Series Of Numbers . . . . .	261
11.3.1 Basic Considerations . . . . .	261
11.3.2 More Tests For Convergence . . . . .	266
11.3.3 Double Series . . . . .	269
11.4 Exercises . . . . .	274
11.5 Taylor Series . . . . .	275
11.5.1 Operations On Power Series . . . . .	277
11.6 Exercises . . . . .	284
11.7 Some Other Theorems . . . . .	286
<b>III Vector Valued Functions</b>	<b>291</b>
<b>12 <math>\mathbb{R}^n</math></b>	<b>293</b>
12.1 Algebra in $\mathbb{R}^n$ . . . . .	294
12.2 Exercises . . . . .	295
12.3 Distance in $\mathbb{R}^n$ . . . . .	296
12.4 Exercises . . . . .	299
12.5 Lines in $\mathbb{R}^n$ . . . . .	300
12.6 Exercises . . . . .	302
12.7 Open And Closed Sets . . . . .	302
12.8 Exercises . . . . .	305
12.9 Vectors . . . . .	305
12.10 Exercises . . . . .	309
<b>13 Vector Products</b>	<b>311</b>
13.1 The Dot Product . . . . .	311
13.2 The Geometric Significance Of The Dot Product . . . . .	313
13.2.1 The Angle Between Two Vectors . . . . .	313
13.2.2 Work And Projections . . . . .	315
13.2.3 The Parabolic Mirror . . . . .	317
13.2.4 The Equation Of A Plane . . . . .	319
13.3 Exercises . . . . .	320
13.4 The Cross Product . . . . .	321
13.4.1 The Distributive Law For The Cross Product . . . . .	324
13.4.2 Torque . . . . .	326
13.4.3 The Box Product . . . . .	328
13.5 Exercises . . . . .	329

13.6	Vector Identities And Notation . . . . .	330
13.7	Exercises . . . . .	332
<b>14</b>	<b>Functions</b>	<b>333</b>
14.1	Exercises . . . . .	334
14.2	Continuous Functions . . . . .	334
14.2.1	Sufficient Conditions For Continuity . . . . .	335
14.3	Exercises . . . . .	335
14.4	Limits Of A Function . . . . .	336
14.5	Exercises . . . . .	339
14.6	The Limit Of A Sequence . . . . .	339
14.6.1	Sequences And Completeness . . . . .	341
14.6.2	Continuity And The Limit Of A Sequence . . . . .	342
14.7	Properties Of Continuous Functions . . . . .	343
14.8	Exercises . . . . .	343
14.9	Some Advanced Calculus . . . . .	344
<b>15</b>	<b>Limits And Derivatives</b>	<b>351</b>
15.1	Limits Of A Vector Valued Function . . . . .	351
15.2	The Derivative And Integral . . . . .	352
15.2.1	Geometric And Physical Significance Of The Derivative . . . . .	354
15.2.2	Differentiation Rules . . . . .	355
15.3	Leibniz's Notation . . . . .	357
15.4	Exercises . . . . .	358
15.5	Newton's Laws Of Motion . . . . .	359
15.5.1	Kinetic Energy . . . . .	363
15.5.2	Impulse And Momentum . . . . .	364
15.6	Exercises . . . . .	365
15.7	Systems Of Ordinary Differential Equations . . . . .	366
15.7.1	Picard Iteration . . . . .	367
15.7.2	Numerical Methods For Differential Equations . . . . .	371
15.8	Exercises . . . . .	373
<b>16</b>	<b>Line Integrals</b>	<b>375</b>
16.1	Arc Length And Orientations . . . . .	375
16.2	Line Integrals And Work . . . . .	378
16.3	Exercises . . . . .	380
16.4	Motion On A Space Curve . . . . .	381
16.5	Exercises . . . . .	384
16.6	Independence Of Parameterization . . . . .	385
<b>17</b>	<b>The Circular Functions Again</b>	<b>387</b>
17.1	The Equations Of Undamped And Damped Oscillation . . . . .	392
17.2	Exercises . . . . .	395
<b>18</b>	<b>Curvilinear Coordinate Systems</b>	<b>397</b>
18.1	Polar Cylindrical And Spherical Coordinates . . . . .	397
18.2	The Acceleration In Polar Coordinates . . . . .	399
18.3	Planetary Motion . . . . .	401
18.4	Exercises . . . . .	406

<b>IV</b>	<b>Functions Of More Than One Variable</b>	<b>409</b>
<b>19</b>	<b>Linear Algebra</b>	<b>411</b>
19.1	Matrices . . . . .	411
19.1.1	Finding The Inverse Of A Matrix . . . . .	417
19.2	Exercises . . . . .	419
19.3	Linear Transformations . . . . .	421
19.3.1	Least Squares Problems . . . . .	423
19.3.2	The Least Squares Regression Line . . . . .	425
19.3.3	The Fredholm Alternative . . . . .	426
19.4	Exercises . . . . .	427
19.5	Moving Coordinate Systems . . . . .	428
19.5.1	The Coriolis Acceleration . . . . .	428
19.5.2	The Coriolis Acceleration On The Rotating Earth . . . . .	432
19.6	Exercises . . . . .	437
19.7	Determinants . . . . .	438
19.8	Exercises . . . . .	445
19.9	The Mathematical Theory Of Determinants . . . . .	447
19.10	Exercises . . . . .	458
19.11	The Determinant And Volume . . . . .	458
19.12	Exercises . . . . .	462
19.13	Linear Systems Of Ordinary Differential Equations . . . . .	462
19.14	Exercises . . . . .	466
<b>20</b>	<b>Functions Of Many Variables</b>	<b>469</b>
20.1	The Graph Of A Function Of Two Variables . . . . .	469
20.2	The Directional Derivative . . . . .	470
20.3	Exercises . . . . .	473
20.4	Mixed Partial Derivatives . . . . .	474
20.5	The Limit Of A Function Of Many Variables . . . . .	475
20.6	Exercises . . . . .	478
20.7	Approximation With A Tangent Plane . . . . .	478
20.8	Exercises . . . . .	484
20.9	Differentiation And The Chain Rule . . . . .	484
20.9.1	The Chain Rule . . . . .	484
20.9.2	Differentiation And The Derivative . . . . .	489
20.10	Exercises . . . . .	492
20.11	The Gradient . . . . .	493
20.12	Exercises . . . . .	495
20.13	Nonlinear Ordinary Differential Equations Local Existence . . . . .	496
20.14	Lagrange Multipliers . . . . .	499
20.15	Exercises . . . . .	504
20.16	Taylor's Formula For Functions Of Many Variables . . . . .	506
20.16.1	Some Linear Algebra . . . . .	508
20.16.2	The Second Derivative Test . . . . .	509
20.17	Exercises . . . . .	513



<b>21 The Riemann Integral On <math>\mathbb{R}^n</math></b>	<b>519</b>
21.1 Methods For Double Integrals . . . . .	519
21.2 Exercises . . . . .	527
21.3 Methods For Triple Integrals . . . . .	528
21.4 Exercises With Answers . . . . .	533
21.5 Exercises . . . . .	537
21.6 Different Coordinates . . . . .	538
21.7 Exercises With Answers . . . . .	543
21.8 Exercises . . . . .	549
21.9 The Moment Of Inertia . . . . .	551
21.9.1 The Spinning Top . . . . .	551
21.9.2 Kinetic Energy . . . . .	554
21.9.3 Finding The Moment Of Inertia And Center Of Mass . . . . .	556
21.10 Exercises . . . . .	557
21.11 Theory Of The Riemann Integral . . . . .	558
21.11.1 Basic Properties . . . . .	560
21.11.2 Iterated Integrals . . . . .	573
21.11.3 Some Observations . . . . .	577
<b>22 The Integral On Other Sets</b>	<b>579</b>
22.1 Exercises With Answers . . . . .	587
22.2 Exercises . . . . .	591
<b>23 Vector Calculus</b>	<b>593</b>
23.1 Divergence And Curl Of A Vector Field . . . . .	593
23.2 Exercises . . . . .	597
23.3 The Divergence Theorem . . . . .	598
23.4 Exercises . . . . .	602
23.5 Some Applications Of The Divergence Theorem . . . . .	603
23.5.1 Hydrostatic Pressure . . . . .	603
23.5.2 Archimedes Law Of Buoyancy . . . . .	603
23.5.3 Equations Of Heat And Diffusion . . . . .	604
23.5.4 Balance Of Mass . . . . .	605
23.5.5 Balance Of Momentum . . . . .	606
23.5.6 The Wave Equation . . . . .	610
23.5.7 A Negative Observation . . . . .	611
23.6 Exercises . . . . .	611
23.7 Stokes Theorem . . . . .	612
23.7.1 Green's Theorem . . . . .	616
23.8 Green's Theorem Again . . . . .	617
23.9 Stoke's Theorem From Green's Theorem . . . . .	619
23.9.1 Conservative Vector Fields . . . . .	622
23.9.2 Maxwell's Equations And The Wave Equation . . . . .	625
23.10 Exercises . . . . .	626
<b>A The Fundamental Theorem Of Algebra</b>	<b>629</b>



# Introduction

Calculus consists of the study of limits of various sorts and the systematic exploitation of the completeness axiom. It was developed by physicists and engineers over a period of several hundred years in order to solve problems from the physical sciences. It is the language by which precision and quantitative predictions for many complicated problems are obtained. It is used to find lengths of curves, areas and volumes of regions which are not bounded by straight lines. It is used to predict and account for the motion of satellites. It is essential in order to solve many maximization problems and it is prerequisite material in order to understand models based on differential equations. These and other applications are discussed to some extent in this book.

It is assumed the reader has a good understanding of algebra on the level of college algebra or what used to be called algebra II along with some exposure to geometry and trigonometry although the book does contain an extensive review of these things. I have tried to keep the book a manageable length in order to focus more on the important ideas. I have also tried to give complete proofs of all theorems in one variable calculus and to at least give plausibility arguments for those in multiple dimensions. Physical models are derived in the usual way through the use of differentials leading to differential equations which are introduced early and used throughout the book as the basis for physical models.

I expect the reader to be able to use a calculator whenever it would be helpful to do so. Many of the exercises will be very troublesome without one. Having said this, calculus is not about using calculators or any other form of technology. I believe that when the syntax and arcane notation associated with technology are presented, these things become the topic of study rather than the concepts of calculus. This is a book on calculus and should not be considered an instruction manual for the use of technology.

Pictures are often helpful in seeing what is going on and there are many pictures in this book for this reason. However, calculus is not about drawing pictures and ultimately rests on logic and definitions. Algebra plays a central role in gaining the sort of understanding which generalizes to higher dimensions where pictures are not available. Therefore, I have emphasized the algebraic aspects of this subject far more than is usual, especially linear algebra which is absolutely essential to understand in order to do multivariable calculus. I have also featured the repeated index summation convention and the usual reduction identities which allow one to discover vector identities.



**Part I**

**Preliminaries**



# The Real Numbers

An understanding of the properties of the real numbers is essential in order to understand calculus. This section contains a review of the algebraic properties of real numbers.

## 2.1 The Number Line And Algebra Of The Real Numbers

To begin with, consider the real numbers, denoted by  $\mathbb{R}$ , as a line extending infinitely far in both directions. In this book, the notation,  $\equiv$  indicates something is being defined. Thus the integers are defined as

$$\mathbb{Z} \equiv \{\cdots -1, 0, 1, \cdots\},$$

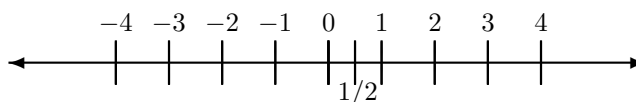
the natural numbers,

$$\mathbb{N} \equiv \{1, 2, \cdots\}$$

and the rational numbers, defined as the numbers which are the quotient of two integers.

$$\mathbb{Q} \equiv \left\{ \frac{m}{n} \text{ such that } m, n \in \mathbb{Z}, n \neq 0 \right\}$$

are each subsets of  $\mathbb{R}$  as indicated in the following picture.



As shown in the picture,  $\frac{1}{2}$  is half way between the number 0 and the number, 1. By analogy, you can see where to place all the other rational numbers. It is assumed that  $\mathbb{R}$  has the following algebra properties, listed here as a collection of assertions called axioms. These properties will not be proved which is why they are called axioms rather than theorems. In general, axioms are statements which are regarded as true. Often these are things which are “self evident” either from experience or from some sort of intuition but this does not have to be the case.

**Axiom 2.1.1**  $x + y = y + x$ , (*commutative law for addition*)

**Axiom 2.1.2**  $x + 0 = x$ , (*additive identity*).

**Axiom 2.1.3** For each  $x \in \mathbb{R}$ , there exists  $-x \in \mathbb{R}$  such that  $x + (-x) = 0$ , (*existence of additive inverse*).

**Axiom 2.1.4**  $(x + y) + z = x + (y + z)$ , (*associative law for addition*).

**Axiom 2.1.5**  $xy = yx$ , (*commutative law for multiplication*).

**Axiom 2.1.6**  $(xy)z = x(yz)$ , (*associative law for multiplication*).

**Axiom 2.1.7**  $1x = x$ , (*multiplicative identity*).

**Axiom 2.1.8** For each  $x \neq 0$ , there exists  $x^{-1}$  such that  $xx^{-1} = 1$ . (*existence of multiplicative inverse*).

**Axiom 2.1.9**  $x(y + z) = xy + xz$ . (*distributive law*).

These axioms are known as the field axioms and any set (there are many others besides  $\mathbb{R}$ ) which has two such operations satisfying the above axioms is called a field. Division and subtraction are defined in the usual way by  $x - y \equiv x + (-y)$  and  $x/y \equiv x(y^{-1})$ . It is assumed that the reader is completely familiar with these axioms in the sense that he or she can do the usual algebraic manipulations taught in high school and junior high algebra courses. The axioms listed above are just a careful statement of exactly what is necessary to make the usual algebraic manipulations valid. A word of advice regarding division and subtraction is in order here. Whenever you feel a little confused about an algebraic expression which involves division or subtraction, think of division as multiplication by the multiplicative inverse as just indicated and think of subtraction as addition of the additive inverse. Thus, when you see  $x/y$ , think  $x(y^{-1})$  and when you see  $x - y$ , think  $x + (-y)$ . In many cases the source of confusion will disappear almost magically. The reason for this is that subtraction and division do not satisfy the associative law. This means there is a natural ambiguity in an expression like  $6 - 3 - 4$ . Do you mean  $(6 - 3) - 4 = -1$  or  $6 - (3 - 4) = 6 - (-1) = 7$ ? It makes a difference doesn't it? However, the so called binary operations of addition and multiplication are associative and so no such confusion will occur. It is conventional to simply do the operations in order of appearance reading from left to right. Thus, if you see  $6 - 3 - 4$ , you would normally interpret it as the first of the above alternatives.

In doing algebra, the following theorem is important and follows from the above axioms. The reasoning which demonstrates this assertion is called a proof. Proofs and definitions are very important in mathematics because they are the means by which "truth" is determined. In mathematics, something is "true" if it follows from axioms using a correct logical argument. Truth is not determined on the basis of experiment or opinions and it is this which makes mathematics useful as a language for describing certain kinds of reality in a precise manner.<sup>1</sup> It is also the definitions and proofs which make the subject of mathematics intellectually worth while. Take these away and it becomes a gray wasteland filled with endless tedium and meaningless manipulations.

In the first part of the following theorem, the claim is made that the additive inverse and the multiplicative inverse are unique. This means that for a given number, only one number has the property that it is an additive inverse and that, given a nonzero number, only one number has the property that it is a multiplicative inverse. The significance of this is that if you are wondering if a given number is the additive inverse of a given number, all you have to do is to check and see if it acts like one.

**Theorem 2.1.10** *The above axioms imply the following.*

1. *The multiplicative inverse and additive inverses are unique.*

---

<sup>1</sup>There are certainly real and important things which should not be described using mathematics because it has nothing to do with these things. For example, feelings and emotions have nothing to do with math.



2.  $0x = 0$ ,  $-(-x) = x$ ,
3.  $(-1)(-1) = 1$ ,  $(-1)x = -x$
4. If  $xy = 0$  then either  $x = 0$  or  $y = 0$ .

**Proof:** Suppose then that  $x$  is a real number and that  $x + y = 0 = x + z$ . It is necessary to verify  $y = z$ . From the above axioms, there exists an additive inverse,  $-x$  for  $x$ . Therefore,

$$-x + 0 = (-x) + (x + y) = (-x) + (x + z)$$

and so by the associative law for addition,

$$((-x) + x) + y = ((-x) + x) + z$$

which implies

$$0 + y = 0 + z.$$

Now by the definition of the additive identity, this implies  $y = z$ . You should prove the multiplicative inverse is unique.

Consider 2. It is desired to verify  $0x = 0$ . From the definition of the additive identity and the distributive law it follows that

$$0x = (0 + 0)x = 0x + 0x.$$

From the existence of the additive inverse and the associative law it follows

$$\begin{aligned} 0 &= (-0x) + 0x = (-0x) + (0x + 0x) \\ &= ((-0x) + 0x) + 0x = 0 + 0x = 0x \end{aligned}$$

To verify the second claim in 2., it suffices to show  $x$  acts like the additive inverse of  $-x$  in order to conclude that  $-(-x) = x$ . This is because it has just been shown that additive inverses are unique. By the definition of additive inverse,

$$x + (-x) = 0$$

and so  $x = -(-x)$  as claimed.

To demonstrate 3.,

$$(-1)(1 + (-1)) = (-1)0 = 0$$

and so using the definition of the multiplicative identity, and the distributive law,

$$(-1) + (-1)(-1) = 0.$$

It follows from 1. and 2. that  $1 = -(-1) = (-1)(-1)$ . To verify  $(-1)x = -x$ , use 2. and the distributive law to write

$$x + (-1)x = x(1 + (-1)) = x0 = 0.$$

Therefore, by the uniqueness of the additive inverse proved in 1., it follows  $(-1)x = -x$  as claimed.

To verify 4., suppose  $x \neq 0$ . Then  $x^{-1}$  exists by the axiom about the existence of multiplicative inverses. Therefore, by 2. and the associative law for multiplication,

$$y = (x^{-1}x)y = x^{-1}(xy) = x^{-1}0 = 0.$$

This proves 4. and completes the proof of this theorem.

Recall the notion of something raised to an integer power. Thus  $y^2 = y \times y$  and  $b^{-3} = \frac{1}{b^3}$  etc.

Also, there are a few conventions related to the order in which operations are performed. Exponents are always done before multiplication. Thus  $xy^2 = x(y^2)$  and is not equal to  $(xy)^2$ . Division or multiplication is always done before addition or subtraction. Thus  $x - y(z + w) = x - [y(z + w)]$  and is not equal to  $(x - y)(z + w)$ . Parentheses are done before anything else. Be very careful of such things since they are a source of mistakes. When you have doubts, insert parentheses to resolve the ambiguities.

Also recall summation notation. If you have not seen this, the following is a short review of this topic.

**Definition 2.1.11** *Let  $x_1, x_2, \dots, x_m$  be numbers. Then*

$$\sum_{j=1}^m x_j \equiv x_1 + x_2 + \dots + x_m.$$

*Thus this symbol,  $\sum_{j=1}^m x_j$  means to take all the numbers,  $x_1, x_2, \dots, x_m$  and add them all up. Note the use of the  $j$  as a generic variable which takes values from 1 up to  $m$ . This notation will be used whenever there are things which can be added, not just numbers.*

As an example of the use of this notation, you should verify the following.

**Example 2.1.12**  $\sum_{k=1}^6 (2k + 1) = 48$ .

Be sure you understand why

$$\sum_{k=1}^{m+1} x_k = \sum_{k=1}^m x_k + x_{m+1}.$$

As a slight generalization of this notation,

$$\sum_{j=k}^m x_j \equiv x_k + \dots + x_m.$$

It is also possible to change the variable of summation.

$$\sum_{j=1}^m x_j = x_1 + x_2 + \dots + x_m$$

while if  $r$  is an integer, the notation requires

$$\sum_{j=1+r}^{m+r} x_{j-r} = x_1 + x_2 + \dots + x_m$$

and so  $\sum_{j=1}^m x_j = \sum_{j=1+r}^{m+r} x_{j-r}$ .

Summation notation will be used throughout the book whenever it is convenient to do so.

Another thing to keep in mind is that you often use letters to represent numbers. Since they represent numbers, you manipulate expressions involving letters in the same manner as you would if they were specific numbers.

**Example 2.1.13** Add the fractions,  $\frac{x}{x^2+y} + \frac{y}{x-1}$ .

You add these just like they were numbers. Write the first expression as  $\frac{x(x-1)}{(x^2+y)(x-1)}$  and the second as  $\frac{y(x^2+y)}{(x-1)(x^2+y)}$ . Then since these have the same common denominator, you add them as follows.

$$\begin{aligned}\frac{x}{x^2+y} + \frac{y}{x-1} &= \frac{x(x-1)}{(x^2+y)(x-1)} + \frac{y(x^2+y)}{(x-1)(x^2+y)} \\ &= \frac{x^2 - x + yx^2 + y^2}{(x^2+y)(x-1)}.\end{aligned}$$

## 2.2 Exercises

1. Consider the expression  $x + y(x + y) - x(y - x) \equiv f(x, y)$ . Find  $f(-1, 2)$ .
2. Show  $-(ab) = (-a)b$ .
3. Show on the number line the effect of adding two positive numbers,  $x$  and  $y$ .
4. Show on the number line the effect of subtracting a positive number from another positive number.
5. Show on the number line the effect of multiplying a number by  $-1$ .
6. Add the fractions  $\frac{x}{x^2-1} + \frac{x-1}{x+1}$ .
7. Find a formula for  $(x+y)^2$ ,  $(x+y)^3$ , and  $(x+y)^4$ . Based on what you observe for these, give a formula for  $(x+y)^8$ .
8. When is it true that  $(x+y)^n = x^n + y^n$ ?
9. Find the error in the following argument. Let  $x = y = 1$ . Then  $xy = y^2$  and so  $xy - x^2 = y^2 - x^2$ . Therefore,  $x(y-x) = (y-x)(y+x)$ . Dividing both sides by  $(y-x)$  yields  $x = x+y$ . Now substituting in what these variables equal yields  $1 = 1+1$ .
10. Find the error in the following argument.  $\sqrt{x^2+1} = x+1$  and so letting  $x = 2$ ,  $\sqrt{5} = 3$ . Therefore,  $5 = 9$ .
11. Find the error in the following. Let  $x = 1$  and  $y = 2$ . Then  $\frac{1}{3} = \frac{1}{x+y} = \frac{1}{x} + \frac{1}{y} = 1 + \frac{1}{2} = \frac{3}{2}$ . Then cross multiplying, yields  $2 = 9$ .
12. Simplify  $\frac{x^2y^4z^{-6}}{x^{-2}y^{-1}z}$ .
13. Simplify the following expressions using correct algebra. In these expressions the variables represent real numbers.
  - (a)  $\frac{x^2y+xy^2+x}{x}$
  - (b)  $\frac{x^2y+xy^2+x}{xy}$
  - (c)  $\frac{x^3+2x^2-x-2}{x+1}$
14. Find the error in the following argument. Let  $x = 3$  and  $y = 1$ . Then  $1 = 3 - 2 = 3 - (3 - 1) = x - y(x - y) = (x - y)(x - y) = 2^2 = 4$ .

15. Verify the following formulas.

$$(a) (x - y)(x + y) = x^2 - y^2$$

$$(b) (x - y)(x^2 + xy + y^2) = x^3 - y^3$$

$$(c) (x + y)(x^2 - xy + y^2) = x^3 + y^3$$

16. Find the error in the following.

$$\frac{xy + y}{x} = y + y = 2y.$$

Now let  $x = 2$  and  $y = 2$  to obtain

$$3 = 4$$

17. Show the rational numbers satisfy the field axioms. You may assume the associative, commutative, and distributive laws hold for the integers.

## 2.3 Order

The real numbers also have an order defined on them. This order can be defined very precisely in terms of a short list of axioms but this will not be done here. Instead, properties which should be familiar are listed here as axioms.

**Definition 2.3.1** *The expression,  $x < y$ , in words, ( $x$  is less than  $y$ ) means  $y$  lies to the right of  $x$  on the number line. The expression  $x > y$ , in words ( $x$  is greater than  $y$ ) means  $x$  is to the right of  $y$  on the number line.  $x \leq y$  if either  $x = y$  or  $x < y$ .  $x \geq y$  if either  $x > y$  or  $x = y$ . A number,  $x$ , is positive if  $x > 0$ .*

If you examine the number line, the following should be fairly reasonable and are listed as axioms, things assumed to be true. I suggest you plug in some numbers to reassure yourself about these axioms.

**Axiom 2.3.2** *The sum of two positive real numbers is positive.*

**Axiom 2.3.3** *The product of two positive real numbers is positive.*

**Axiom 2.3.4** *For a given real number  $x$ , one and only one of the following alternatives holds. Either  $x$  is positive,  $x = 0$ , or  $-x$  is positive.*

**Axiom 2.3.5** *If  $x < y$  and  $y < z$  then  $x < z$  (Transitive law).*

**Axiom 2.3.6** *If  $x < y$  then  $x + z < y + z$  (addition to an inequality).*

**Axiom 2.3.7** *If  $x \leq 0$  and  $y \leq 0$ , then  $xy \geq 0$ .*

**Axiom 2.3.8** *If  $x > 0$  then  $x^{-1} > 0$ .*

**Axiom 2.3.9** *If  $x < 0$  then  $x^{-1} < 0$ .*

**Axiom 2.3.10** *If  $x < y$  then  $xz < yz$  if  $z > 0$ , (multiplication of an inequality by a positive number).*

**Axiom 2.3.11** *If  $x < y$  and  $z < 0$ , then  $xz > yz$  (multiplication of an inequality by a negative number).*

**Axiom 2.3.12** Each of the above holds with  $>$  and  $<$  replaced by  $\geq$  and  $\leq$  respectively except for 2.3.8 and 2.3.9 in which it is also necessary to stipulate that  $x \neq 0$ .

**Axiom 2.3.13** For any  $x$  and  $y$ , exactly one of the following must hold. Either  $x = y$ ,  $x < y$ , or  $x > y$  (trichotomy).

Note that trichotomy could be stated by saying  $x \leq y$  or  $y \leq x$ .

**Example 2.3.14** Solve the inequality  $2x + 4 \leq x - 8$

Subtract  $2x$  from both sides to yield  $4 \leq -x - 8$ . Next add 8 to both sides to get  $12 \leq -x$ . Then multiply both sides by  $(-1)$  to obtain  $x \leq -12$ . Alternatively, subtract  $x$  from both sides to get  $x + 4 \leq -8$ . Then subtract 4 from both sides to obtain  $x \leq -12$ .

**Example 2.3.15** Solve the inequality  $(x + 1)(2x - 3) \geq 0$ .

If this is to hold, either both of the factors,  $x + 1$  and  $2x - 3$  are nonnegative or they are both nonpositive. The first case yields  $x + 1 \geq 0$  and  $2x - 3 \geq 0$  so  $x \geq -1$  and  $x \geq \frac{3}{2}$  yielding  $x \geq \frac{3}{2}$ . The second case yields  $x + 1 \leq 0$  and  $2x - 3 \leq 0$  which implies  $x \leq -1$  and  $x \leq \frac{3}{2}$ . Therefore, the solution to this inequality is  $x \leq -1$  or  $x \geq \frac{3}{2}$ .

**Example 2.3.16** Solve the inequality  $(x)(x + 2) \geq -4$

Here the problem is to find  $x$  such that  $x^2 + 2x + 4 \geq 0$ . However,  $x^2 + 2x + 4 = (x + 1)^2 + 3 \geq 0$  for all  $x$ . Therefore, the solution to this problem is all  $x \in \mathbb{R}$ .

To simplify the way such things are written, involves set notation. This is described next.

### 2.3.1 Set Notation

A set is just a collection of things called elements. For example  $\{1, 2, 3, 8\}$  would be a set consisting of the elements 1, 2, 3, and 8. To indicate that 3 is an element of  $\{1, 2, 3, 8\}$ , it is customary to write  $3 \in \{1, 2, 3, 8\}$ .  $9 \notin \{1, 2, 3, 8\}$  means 9 is not an element of  $\{1, 2, 3, 8\}$ . Sometimes a rule specifies a set. For example you could specify a set as all integers larger than 2. This would be written as  $S = \{x \in \mathbb{Z} : x > 2\}$ . This notation says: the set of all integers,  $x$ , such that  $x > 2$ .

If  $A$  and  $B$  are sets with the property that every element of  $A$  is an element of  $B$ , then  $A$  is a subset of  $B$ . For example,  $\{1, 2, 3, 8\}$  is a subset of  $\{1, 2, 3, 4, 5, 8\}$ , in symbols,  $\{1, 2, 3, 8\} \subseteq \{1, 2, 3, 4, 5, 8\}$ . The same statement about the two sets may also be written as  $\{1, 2, 3, 4, 5, 8\} \supseteq \{1, 2, 3, 8\}$ .

The union of two sets is the set consisting of everything which is contained in at least one of the sets,  $A$  or  $B$ . As an example of the union of two sets,  $\{1, 2, 3, 8\} \cup \{3, 4, 7, 8\} = \{1, 2, 3, 4, 7, 8\}$  because these numbers are those which are in at least one of the two sets. In general

$$A \cup B \equiv \{x : x \in A \text{ or } x \in B\}.$$

Be sure you understand that something which is in both  $A$  and  $B$  is in the union. It is not an exclusive or.

The intersection of two sets,  $A$  and  $B$  consists of everything which is in both of the sets. Thus  $\{1, 2, 3, 8\} \cap \{3, 4, 7, 8\} = \{3, 8\}$  because 3 and 8 are those elements the two sets have in common. In general,

$$A \cap B \equiv \{x : x \in A \text{ and } x \in B\}.$$

When with real numbers,  $[a, b]$  denotes the set of real numbers,  $x$ , such that  $a \leq x \leq b$  and  $[a, b)$  denotes the set of real numbers such that  $a \leq x < b$ .  $(a, b)$  consists of the set of real numbers,  $x$  such that  $a < x < b$  and  $(a, b]$  indicates the set of numbers,  $x$  such that  $a < x \leq b$ .  $[a, \infty)$  means the set of all numbers,  $x$  such that  $x \geq a$  and  $(-\infty, a]$  means the set of all real numbers which are less than or equal to  $a$ . These sorts of sets of real numbers are called intervals. The two points,  $a$  and  $b$  are called endpoints of the interval. Other intervals such as  $(-\infty, b)$  are defined by analogy to what was just explained. In general, the curved parenthesis indicates the end point it sits next to is not included while the square parenthesis indicates this end point is included. The reason that there will always be a curved parenthesis next to  $\infty$  or  $-\infty$  is that these are not real numbers. Therefore, they cannot be included in any set of real numbers.

A special set which needs to be given a name is the empty set also called the null set, denoted by  $\emptyset$ . Thus  $\emptyset$  is defined as the set which has no elements in it. Mathematicians like to say the empty set is a subset of every set. The reason they say this is that if it were not so, there would have to exist a set,  $A$ , such that  $\emptyset$  has something in it which is not in  $A$ . However,  $\emptyset$  has nothing in it and so the least intellectual discomfort is achieved by saying  $\emptyset \subseteq A$ .

If  $A$  and  $B$  are two sets,  $A \setminus B$  denotes the set of things which are in  $A$  but not in  $B$ . Thus

$$A \setminus B \equiv \{x \in A : x \notin B\}.$$

Set notation is used whenever convenient.

To illustrate the use of this notation consider the same three examples of inequalities.

**Example 2.3.17** Solve the inequality  $2x + 4 \leq x - 8$

This was worked earlier and  $x \leq -12$  was the answer. This is written as  $(-\infty, -12]$ .

**Example 2.3.18** Solve the inequality  $(x + 1)(2x - 3) \geq 0$ .

This was worked earlier and  $x \leq -1$  or  $x \geq \frac{3}{2}$  was the answer. In terms of set notation this is denoted by  $(-\infty, -1] \cup [\frac{3}{2}, \infty)$ .

**Example 2.3.19** Solve the inequality  $(x)(x + 2) \geq -4$

Recall this inequality was true for any value of  $x$ . It is written as  $\mathbb{R}$  or  $(-\infty, \infty)$ .

## 2.4 Exercises With Answers

1. Solve  $(3x + 1)(x - 2) \leq 0$ .

This happens when the two factors have different signs. Thus either  $3x + 1 \leq 0$  and  $x - 2 \geq 0$  in which case  $x \leq -\frac{1}{3}$  and  $x \geq 2$ , a situation which never occurs, or else  $3x + 1 \geq 0$  and  $x - 2 \leq 0$  so  $x \geq -\frac{1}{3}$  and  $x \leq 2$ . Written as  $[-\frac{1}{3}, 2]$ .

2. Solve  $(3x + 1)(x - 2) > 0$ .

This is just everything not included in the above problem. Thus the answer would be  $(-\infty, -\frac{1}{3}) \cup (2, \infty)$ .

3. Solve  $\frac{x+1}{2x-2} < 0$ .

Note that  $\frac{x+1}{2x-2}$  is positive if  $x > 1$ , negative if  $x \in (-1, 1)$ , and nonnegative if  $x \leq -1$ . Therefore, the answer is  $(-1, 1)$ . To identify the interesting intervals, all that was necessary to do was to look at the two factors,  $(x + 1)$  and  $(2x - 2)$  and determine where these equal zero.

4. Solve  $\frac{3x+7}{x^2+2x+1} \geq 1$ .

On something like this, subtract 1 from both sides to get

$$\frac{6+x-x^2}{x^2+2x+1} = \frac{(3-x)(2+x)}{(x+1)^2}.$$

When  $x = 3$  or  $x = -2$ , this equals zero. For  $x \in (-2, 3)$  the expression is positive and it is negative if  $x > 3$  or if  $x < -2$ . Therefore, the answer is  $[-2, 3]$ .

## 2.5 Exercises

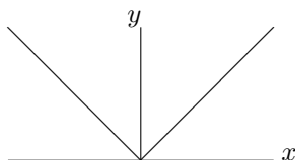
1. Solve  $(3x+2)(x-3) \leq 0$ .
2. Solve  $(3x+2)(x-3) > 0$ .
3. Solve  $\frac{x+2}{3x-2} < 0$ .
4. Solve  $\frac{x+1}{x+3} < 1$ .
5. Solve  $(x-1)(2x+1) \leq 2$ .
6. Solve  $(x-1)(2x+1) > 2$ .
7. Solve  $x^2 - 2x \leq 0$ .
8. Solve  $(x+2)(x-2)^2 \leq 0$ .
9. Solve  $\frac{3x-4}{x^2+2x+2} \geq 0$ .
10. Solve  $\frac{3x+9}{x^2+2x+1} \geq 1$ .
11. Solve  $\frac{x^2+2x+1}{3x+7} < 1$ .

## 2.6 The Absolute Value

A fundamental idea is the absolute value of a number. This is important because the absolute value defines distance on  $\mathbb{R}$ . How far away from 0 is the number 3? How about the number  $-3$ ? Look at the number line and observe they are both 3 units away from 0. To describe this algebraically,

**Definition 2.6.1**  $|x| \equiv \begin{cases} x & \text{if } x \geq 0, \\ -x & \text{if } x < 0. \end{cases}$

Thus  $|x|$  can be thought of as the distance between  $x$  and 0. It may be useful to think of this function in terms of its graph if you recall the notion of the graph of a function.



The following is a fundamental theorem about the absolute value.

**Theorem 2.6.2**  $|xy| = |x| |y|$ .

**Proof:** If both  $x, y \leq 0$ , then  $|xy| = xy$  because in this case  $xy \geq 0$  while

$$|x||y| = (-x)(-y) = (-1)x(-1)y = (-1)(-1)xy = xy.$$

Therefore, in this case the result of the theorem is verified. You should verify the other cases, both  $x, y \geq 0$  and  $x \leq 0$  while  $y \geq 0$ .

This theorem is the basis for the following fundamental result which is of major importance in calculus.

**Theorem 2.6.3** *The following inequalities hold.*

$$|x + y| \leq |x| + |y|, \quad ||x| - |y|| \leq |x - y|.$$

*Either of these inequalities may be called the triangle inequality.*

**Proof:** By Theorem 2.6.2,

$$\begin{aligned} |x + y|^2 &= |(x + y)^2| = (x + y)^2 \\ &= x^2 + y^2 + 2xy \leq x^2 + y^2 + 2|x||y| \\ &= |x|^2 + |y|^2 + 2|x||y| = (|x| + |y|)^2. \end{aligned}$$

Now note that if  $0 \leq a \leq b$  then  $0 \leq a^2 \leq ab \leq b^2$  and that if  $a, b \geq 0$  then if  $a^2 \leq b^2$  it follows that  $b^2 \geq ba \geq a^2$  and so  $b \geq a$  (see the above axioms. Multiply by  $a^{-1}$  if  $a \neq 0$ .) Applying this observation to the above inequality,

$$|x + y| \leq |x| + |y|.$$

This verifies the first of these inequalities. To obtain the second one, note

$$\begin{aligned} |x| &= |x - y + y| \\ &\leq |x - y| + |y| \end{aligned}$$

and so

$$|x| - |y| \leq |x - y| \tag{2.1}$$

Now switch the letters to obtain

$$|y| - |x| \leq |y - x| = |x - y|. \tag{2.2}$$

Therefore,

$$||x| - |y|| \leq |x - y|$$

because if  $|x| - |y| \geq 0$ , then the conclusion follows from (2.1) while if  $|x| - |y| \leq 0$ , the conclusion follows from (2.2). This proves the theorem.

Note there is an inequality involved. Consider the following.

$$|3 + (-2)| = |1| = 1$$

while

$$|3| + |(-2)| = 3 + 2 = 5.$$

You observe that  $5 > 1$  and so it is important to remember that the triangle inequality is an inequality.

**Example 2.6.4** *Solve the equation  $|x - 1| = 2$*



This will be true when  $x - 1 = 2$  or when  $x - 1 = -2$ . Therefore, there are two solutions to this problem,  $x = 3$  or  $x = -1$ .

**Example 2.6.5** Solve the inequality  $|2x - 1| < 2$

From the number line, it is necessary to have  $2x - 1$  between  $-2$  and  $2$  because the inequality says that the distance from  $2x - 1$  to  $0$  is less than  $2$ . Therefore,  $-2 < 2x - 1 < 2$  and so  $-1/2 < x < 3/2$ . In other words,  $-1/2 < x$  and  $x < 3/2$ .

**Example 2.6.6** Solve the inequality  $|2x - 1| > 2$ .

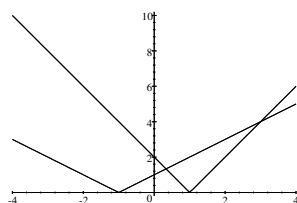
This happens if  $2x - 1 > 2$  or if  $2x - 1 < -2$ . Thus the solution is  $x > 3/2$  or  $x < -1/2$ ,  $(\frac{3}{2}, \infty) \cup (-\infty, -\frac{1}{2})$ .

**Example 2.6.7** Solve  $|x + 1| = |2x - 2|$

There are two ways this can happen. It could be the case that  $x + 1 = 2x - 2$  in which case  $x = 3$  or alternatively,  $x + 1 = 2 - 2x$  in which case  $x = 1/3$ .

**Example 2.6.8** Solve  $|x + 1| \leq |2x - 2|$

In order to keep track of what is happening, it is a very good idea to graph the two relations,  $y = |x + 1|$  and  $y = |2x - 2|$  on the same set of coordinate axes. This is not a hard job.  $|x + 1| = x + 1$  when  $x > -1$  and  $|x + 1| = -1 - x$  when  $x \leq -1$ . Therefore, it is not hard to draw its graph. Similar considerations apply to the other relation. The result is



Equality holds exactly when  $x = 3$  or  $x = \frac{1}{3}$  as in the preceding example. Consider  $x$  between  $\frac{1}{3}$  and  $3$ . You can see these values of  $x$  do not solve the inequality. For example  $x = 1$  does not work. Therefore,  $(\frac{1}{3}, 3)$  must be excluded. The values of  $x$  larger than  $3$  do not produce equality so either  $|x + 1| < |2x - 2|$  for these points or  $|2x - 2| < |x + 1|$  for these points. Checking examples, you see the first of the two cases is the one which holds. Therefore,  $[3, \infty)$  is included. Similar reasoning obtains  $(-\infty, \frac{1}{3}]$ . It follows the solution set to this inequality is  $(-\infty, \frac{1}{3}] \cup [3, \infty)$ .

**Example 2.6.9** Obtain a number,  $\delta$ , such that if  $|x - 2| < \delta$ , then  $|x^2 - 4| < 1/10$ .

If  $|x - 2| < 1$ , then  $||x| - 2| < 1$  and so  $|x| < 3$ . Therefore, if  $|x - 2| < 1$ ,

$$\begin{aligned} |x^2 - 4| &= |x + 2| |x - 2| \\ &\leq (|x| + 2) |x - 2| \\ &\leq 5 |x - 2|. \end{aligned}$$

Therefore, if  $|x - 2| < \frac{1}{50}$ , the desired inequality will hold. Note that some of this is arbitrary. For example, if  $|x - 2| < 3$ , then  $||x| - 2| < 3$  and so  $|x| < 5$ . Therefore, for such  $x$ ,

$$\begin{aligned} |x^2 - 4| &= |x + 2| |x - 2| \\ &\leq (|x| + 2) |x - 2| \\ &\leq 7 |x - 2|. \end{aligned}$$

and so it would also suffice to take  $|x - 2| < \frac{1}{70}$ . The example is about the existence of a number which has a certain property, not the question of finding a particular such number. There are infinitely many which will work because if you have found one, then any which is smaller will also work.

**Example 2.6.10** Suppose  $\varepsilon > 0$  is a given positive number. Obtain a number,  $\delta > 0$ , such that if  $|x - 1| < \delta$ , then  $|x^2 - 1| < \varepsilon$ .

First of all, note  $|x^2 - 1| = |x - 1||x + 1| \leq (|x| + 1)|x - 1|$ . Now if  $|x - 1| < 1$ , it follows  $|x| < 2$  and so for  $|x - 1| < 1$ ,

$$|x^2 - 1| < 3|x - 1|.$$

Now let  $\delta = \min(1, \frac{\varepsilon}{3})$ . This notation means to take the minimum of the two numbers, 1 and  $\frac{\varepsilon}{3}$ . Then if  $|x - 1| < \delta$ ,

$$|x^2 - 1| < 3|x - 1| < 3\frac{\varepsilon}{3} = \varepsilon.$$

## 2.7 Exercises

1. Solve  $|x + 1| = |2x - 3|$ .
2. Solve  $|3x + 1| < 8$ . Give your answer in terms of intervals on the real line.
3. Solve  $|x + 2| < |3x - 3|$ .
4. Tell when equality holds in the triangle inequality.
5. Solve  $|x + 2| \leq 8 + |2x - 4|$ .
6. Verify the axioms for order listed above are reasonable by consideration of the number line. In particular, show that if  $x \leq z$  and  $y < 0$  then  $xy \geq yz$ .
7. Solve  $(x + 1)(2x - 2)x \geq 0$ .
8. Solve  $\frac{x+3}{2x+1} > 1$ .
9. Solve  $\frac{x+2}{3x+1} > 2$ .
10. Describe the set of numbers,  $a$  such that there is no solution to  $|x + 1| = 4 - |x + a|$ .
11. Suppose  $0 < a < b$ . Show  $a^{-1} > b^{-1}$ .
12. Show that if  $|x - 6| < 1$ , then  $|x| < 7$ .
13. Suppose  $|x - 8| < 2$ . How large can  $|x - 5|$  be?
14. Obtain a number,  $\delta > 0$ , such that if  $|x - 1| < \delta$ , then  $|x^2 - 1| < 1/10$ .
15. Obtain a number,  $\delta > 0$ , such that if  $|x - 4| < \delta$ , then  $|\sqrt{x} - 2| < 1/10$ .
16. Suppose  $\varepsilon > 0$  is a given positive number. Obtain a number,  $\delta > 0$ , such that if  $|x - 1| < \delta$ , then  $|\sqrt{x} - 1| < \varepsilon$ . **Hint:** This  $\delta$  will depend in some way on  $\varepsilon$ . You need to tell how.

## 2.8 Well Ordering Principle And Archimedian Property

**Definition 2.8.1** A set is well ordered if every nonempty subset  $S$ , contains a smallest element  $z$  having the property that  $z \leq x$  for all  $x \in S$ .

**Axiom 2.8.2** Any set of integers larger than a given number is well ordered.

In particular, the natural numbers defined as

$$\mathbb{N} \equiv \{1, 2, \dots\}$$

is well ordered.

The above axiom implies the principle of mathematical induction.

**Theorem 2.8.3** (Mathematical induction) A set  $S \subseteq \mathbb{Z}$ , having the property that  $a \in S$  and  $n + 1 \in S$  whenever  $n \in S$  contains all integers  $x \in \mathbb{Z}$  such that  $x \geq a$ .

**Proof:** Let  $T \equiv ([a, \infty) \cap \mathbb{Z}) \setminus S$ . Thus  $T$  consists of all integers larger than or equal to  $a$  which are not in  $S$ . The theorem will be proved if  $T = \emptyset$ . If  $T \neq \emptyset$  then by the well ordering principle, there would have to exist a smallest element of  $T$ , denoted as  $b$ . It must be the case that  $b > a$  since by definition,  $a \notin T$ . Then the integer,  $b - 1 \geq a$  and  $b - 1 \notin S$  because if  $b \in S$ , then  $b - 1 + 1 = b \in S$  by the assumed property of  $S$ . Therefore,  $b - 1 \in ([a, \infty) \cap \mathbb{Z}) \setminus S = T$  which contradicts the choice of  $b$  as the smallest element of  $T$ . ( $b - 1$  is smaller.) Since a contradiction is obtained by assuming  $T \neq \emptyset$ , it must be the case that  $T = \emptyset$  and this says that everything in  $[a, \infty) \cap \mathbb{Z}$  is also in  $S$ .

Mathematical induction is a very useful device for proving theorems about the integers.

**Example 2.8.4** Prove by induction that  $\sum_{k=1}^n k^2 = \frac{n(n+1)(2n+1)}{6}$ .

By inspection, if  $n = 1$  then the formula is true. The sum yields 1 and so does the formula on the right. Suppose this formula is valid for some  $n \geq 1$  where  $n$  is an integer. Then

$$\begin{aligned} \sum_{k=1}^{n+1} k^2 &= \sum_{k=1}^n k^2 + (n+1)^2 \\ &= \frac{n(n+1)(2n+1)}{6} + (n+1)^2. \end{aligned}$$

The step going from the first to the second line is based on the assumption that the formula is true for  $n$ . This is called the induction hypothesis. Now simplify the expression in the second line,

$$\frac{n(n+1)(2n+1)}{6} + (n+1)^2.$$

This equals

$$(n+1) \left( \frac{n(2n+1)}{6} + (n+1) \right)$$

and

$$\begin{aligned} \frac{n(2n+1)}{6} + (n+1) &= \frac{6(n+1) + 2n^2 + n}{6} \\ &= \frac{(n+2)(2n+3)}{6} \end{aligned}$$

Therefore,

$$\begin{aligned}\sum_{k=1}^{n+1} k^2 &= \frac{(n+1)(n+2)(2n+3)}{6} \\ &= \frac{(n+1)((n+1)+1)(2(n+1)+1)}{6},\end{aligned}$$

showing the formula holds for  $n+1$  whenever it holds for  $n$ . This proves the formula by mathematical induction.

**Example 2.8.5** Show that for all  $n \in \mathbb{N}$ ,  $\frac{1}{2} \cdot \frac{3}{4} \cdots \frac{2n-1}{2n} < \frac{1}{\sqrt{2n+1}}$ .

If  $n = 1$  this reduces to the statement that  $\frac{1}{2} < \frac{1}{\sqrt{3}}$  which is obviously true. Suppose then that the inequality holds for  $n$ . Then

$$\begin{aligned}\frac{1}{2} \cdot \frac{3}{4} \cdots \frac{2n-1}{2n} \cdot \frac{2n+1}{2n+2} &< \frac{1}{\sqrt{2n+1}} \cdot \frac{2n+1}{2n+2} \\ &= \frac{\sqrt{2n+1}}{2n+2}.\end{aligned}$$

The theorem will be proved if this last expression is less than  $\frac{1}{\sqrt{2n+3}}$ . This happens if and only if

$$\left( \frac{1}{\sqrt{2n+3}} \right)^2 = \frac{1}{2n+3} > \frac{2n+1}{(2n+2)^2}$$

which occurs if and only if  $(2n+2)^2 > (2n+3)(2n+1)$  and this is clearly true which may be seen from expanding both sides. This proves the inequality.

Lets review the process just used. If  $S$  is the set of integers at least as large as 1 for which the formula holds, the first step was to show  $1 \in S$  and then that whenever  $n \in S$ , it follows  $n+1 \in S$ . Therefore, by the principle of mathematical induction,  $S$  contains  $[1, \infty) \cap \mathbb{Z}$ , all positive integers. In doing an inductive proof of this sort, the set,  $S$  is normally not mentioned. One just verifies the steps above. First show the thing is true for some  $a \in \mathbb{Z}$  and then verify that whenever it is true for  $m$  it follows it is also true for  $m+1$ . When this has been done, the theorem has been proved for all  $m \geq a$ .

**Definition 2.8.6** The Archimedian property states that whenever  $x \in \mathbb{R}$ , and  $a > 0$ , there exists  $n \in \mathbb{N}$  such that  $na > x$ .

**Axiom 2.8.7**  $\mathbb{R}$  has the Archimedian property.

This is not hard to believe. Just look at the number line. This Archimedian property is quite important because it shows every real number is smaller than some integer. It also can be used to verify a very important property of the rational numbers.

**Theorem 2.8.8** Suppose  $x < y$  and  $y - x > 1$ . Then there exists an integer,  $l \in \mathbb{Z}$ , such that  $x < l < y$ . If  $x$  is an integer, there is no integer  $y$  satisfying  $x < y < x+1$ .

**Proof:** Let  $x$  be the smallest positive integer. Not surprisingly,  $x = 1$  but this can be proved. If  $x < 1$  then  $x^2 < x$  contradicting the assertion that  $x$  is the smallest natural number. Therefore, 1 is the smallest natural number. This shows there is no integer,  $y$ , satisfying  $x < y < x+1$  since otherwise, you could subtract  $x$  and conclude  $0 < y - x < 1$  for some integer  $y - x$ .

Now suppose  $y - x > 1$  and let

$$S \equiv \{w \in \mathbb{N} : w \geq y\}.$$

The set  $S$  is nonempty by the Archimedian property. Let  $k$  be the smallest element of  $S$ . Therefore,  $k - 1 < y$ . Either  $k - 1 \leq x$  or  $k - 1 > x$ . If  $k - 1 \leq x$ , then

$$y - x \leq y - (k - 1) = \overbrace{y - k}^{\leq 0} + 1 \leq 1$$

contrary to the assumption that  $y - x > 1$ . Therefore,  $x < k - 1 < y$  and this proves the theorem with  $l = k - 1$ .

It is the next theorem which gives the density of the rational numbers. This means that for any real number, there exists a rational number arbitrarily close to it.

**Theorem 2.8.9** *If  $x < y$  then there exists a rational number  $r$  such that  $x < r < y$ .*

**Proof:** Let  $n \in \mathbb{N}$  be large enough that

$$n(y - x) > 1.$$

Thus  $(y - x)$  added to itself  $n$  times is larger than 1. Thus,

$$n(y - x) = ny + n(-x) = ny - nx > 1.$$

It follows from Theorem 2.8.8 there exists  $m \in \mathbb{Z}$  such that

$$nx < m < ny$$

and so take  $r = m/n$ .

**Definition 2.8.10** *A set,  $S \subseteq \mathbb{R}$  is dense in  $\mathbb{R}$  if whenever  $a < b$ ,  $S \cap (a, b) \neq \emptyset$ .*

Thus the above theorem says  $\mathbb{Q}$  is “dense” in  $\mathbb{R}$ .

You probably saw the process of division in elementary school. Even though you saw it at a young age it is very profound and quite difficult to understand. Suppose you want to do the following problem  $\frac{79}{22}$ . What did you do? You likely did a process of long division which gave the following result.

$$\frac{79}{22} = 3 \text{ with remainder } 13.$$

This meant

$$79 = 3(22) + 13.$$

You were given two numbers, 79 and 22 and you wrote the first as some multiple of the second added to a third number which was smaller than the second number. Can this always be done? The answer is in the next theorem and depends here on the Archimedian property of the real numbers.

**Theorem 2.8.11** *Suppose  $0 < a$  and let  $b \geq 0$ . Then there exists a unique integer  $p$  and real number  $r$  such that  $0 \leq r < a$  and  $b = pa + r$ .*

**Proof:** Let  $S \equiv \{n \in \mathbb{N} : an > b\}$ . By the Archimedean property this set is nonempty. Let  $p + 1$  be the smallest element of  $S$ . Then  $pa \leq b$  because  $p + 1$  is the smallest in  $S$ . Therefore,

$$r \equiv b - pa \geq 0.$$

If  $r \geq a$  then  $b - pa \geq a$  and so  $b \geq (p + 1)a$  contradicting  $p + 1 \in S$ . Therefore,  $r < a$  as desired.

To verify uniqueness of  $p$  and  $r$ , suppose  $p_i$  and  $r_i$ ,  $i = 1, 2$ , both work and  $r_2 > r_1$ . Then a little algebra shows

$$p_1 - p_2 = \frac{r_2 - r_1}{a} \in (0, 1).$$

Thus  $p_1 - p_2$  is an integer between 0 and 1, contradicting Theorem 2.8.8. The case that  $r_1 > r_2$  cannot occur either by similar reasoning. Thus  $r_1 = r_2$  and it follows that  $p_1 = p_2$ .

This theorem is called the Euclidean algorithm when  $a$  and  $b$  are integers.

## 2.9 Exercises

1. The Archimedean property implies the rational numbers are dense in  $\mathbb{R}$ . Now consider the numbers which are of the form  $\frac{k}{2^m}$  where  $k \in \mathbb{Z}$  and  $m \in \mathbb{N}$ . Using the number line, demonstrate that the numbers of this form are also dense in  $\mathbb{R}$ .
2. Show there is no smallest number in  $(0, 1)$ . Recall  $(0, 1)$  means the real numbers which are strictly larger than 0 and smaller than 1.
3. Show there is no smallest number in  $\mathbb{Q} \cap (0, 1)$ .
4. Show that if  $S \subseteq \mathbb{R}$  and  $S$  is well ordered with respect to the usual order on  $\mathbb{R}$  then  $S$  cannot be dense in  $\mathbb{R}$ .
5. Prove by induction that  $\sum_{k=1}^n k^3 = \frac{1}{4}n^4 + \frac{1}{2}n^3 + \frac{1}{4}n^2$ .
6. It is a fine thing to be able to prove a theorem by induction but it is even better to be able to come up with a theorem to prove in the first place. Derive a formula for  $\sum_{k=1}^n k^4$  in the following way. Look for a formula in the form  $An^5 + Bn^4 + Cn^3 + Dn^2 + En + F$ . Then try to find the constants  $A, B, C, D, E$ , and  $F$  such that things work out right. In doing this, show

$$\begin{aligned} (n+1)^4 = & \\ \left( A(n+1)^5 + B(n+1)^4 + C(n+1)^3 + D(n+1)^2 + E(n+1) + F \right) & \\ - An^5 + Bn^4 + Cn^3 + Dn^2 + En + F & \end{aligned}$$

and so some progress can be made by matching the coefficients. When you get your answer, prove it is valid by induction.

7. Prove by induction that whenever  $n \geq 2$ ,  $\sum_{k=1}^n \frac{1}{\sqrt{k}} > \sqrt{n}$ .
8. If  $r \neq 0$ , show by induction that  $\sum_{k=1}^n ar^k = a \frac{r^{n+1}}{r-1} - a \frac{r}{r-1}$ .
9. Prove by induction that  $\sum_{k=1}^n k = \frac{n(n+1)}{2}$ .
10. Let  $a$  and  $d$  be real numbers. Find a formula for  $\sum_{k=1}^n (a + kd)$  and then prove your result by induction.

11. Consider the geometric series,  $\sum_{k=1}^n ar^{k-1}$ . Prove by induction that if  $r \neq 1$ , then

$$\sum_{k=1}^n ar^{k-1} = \frac{a - ar^n}{1 - r}.$$

12. This problem is a continuation of Problem 11. You put money in the bank and it accrues interest at the rate of  $r$  per payment period. These terms need a little explanation. If the payment period is one month, and you started with \$100 then the amount at the end of one month would equal  $100(1+r) = 100 + 100r$ . In this the second term is the interest and the first is called the principal. Now you have  $100(1+r)$  in the bank. How much will you have at the end of the second month? By analogy to what was just done it would equal

$$100(1+r) + 100(1+r)r = 100(1+r)^2.$$

In general, the amount you would have at the end of  $n$  months would be  $100(1+r)^n$ . (When a bank says they offer 6% compounded monthly, this means  $r$ , the rate per payment period equals .06/12.) In general, suppose you start with  $P$  and it sits in the bank for  $n$  payment periods. Then at the end of the  $n^{\text{th}}$  payment period, you would have  $P(1+r)^n$  in the bank. In an ordinary annuity, you make payments,  $P$  at the end of each payment period, the first payment at the end of the first payment period. Thus there are  $n$  payments in all. Each accrue interest at the rate of  $r$  per payment period. Using Problem 11, find a formula for the amount you will have in the bank at the end of  $n$  payment periods? This is called the future value of an ordinary annuity. **Hint:** The first payment sits in the bank for  $n-1$  payment periods and so this payment becomes  $P(1+r)^{n-1}$ . The second sits in the bank for  $n-2$  payment periods so it grows to  $P(1+r)^{n-2}$ , etc.

13. Now suppose you want to buy a house by making  $n$  equal monthly payments. Typically,  $n$  is pretty large, 360 for a thirty year loan. Clearly a payment made 10 years from now can't be considered as valuable to the bank as one made today. This is because the one made today could be invested by the bank and having accrued interest for 10 years would be far larger. So what is a payment made at the end of  $k$  payment periods worth today assuming money is worth  $r$  per payment period? Shouldn't it be the amount,  $Q$  which when invested at a rate of  $r$  per payment period would yield  $P$  at the end of  $k$  payment periods? Thus from Problem 12  $Q(1+r)^k = P$  and so  $Q = P(1+r)^{-k}$ . Thus this payment of  $P$  at the end of  $n$  payment periods, is worth  $P(1+r)^{-k}$  to the bank right now. It follows the amount of the loan should equal the sum of these "discounted payments". That is, letting  $A$  be the amount of the loan,

$$A = \sum_{k=1}^n P(1+r)^{-k}.$$

Using Problem 11, find a formula for the right side of the above formula. This is called the present value of an ordinary annuity.

14. Suppose the available interest rate is 7% per year and you want to take a loan for \$100,000 with the first monthly payment at the end of the first month. If you want to pay off the loan in 20 years, what should the monthly payments be? **Hint:** The rate per payment period is .07/12. See the formula you got in Problem 13 and solve for  $P$ .

15. Consider the first five rows of Pascal's<sup>2</sup> triangle

$$\begin{array}{c} 1 \\ 1 \ 1 \\ 1 \ 2 \ 1 \\ 1 \ 3 \ 3 \ 1 \\ 1 \ 4 \ 6 \ 4 \ 1 \end{array}$$

What would the sixth row be? Now consider that  $(x+y)^1 = 1x + 1y$ ,  $(x+y)^2 = x^2 + 2xy + y^2$ , and  $(x+y)^3 = x^3 + 3x^2y + 3xy^2 + y^3$ . Give a conjecture about that  $(x+y)^5$  would be.

16. Based on Problem 15 conjecture a formula for  $(x+y)^n$  and prove your conjecture by induction. **Hint:** Letting the numbers of the  $n^{\text{th}}$  row of Pascal's triangle be denoted by  $\binom{n}{0}, \binom{n}{1}, \dots, \binom{n}{n}$  in reading from left to right, there is a relation between the numbers on the  $(n+1)^{\text{st}}$  row and those on the  $n^{\text{th}}$  row, the relation being  $\binom{n+1}{k} = \binom{n}{k} + \binom{n}{k-1}$ . This is used in the inductive step.
17. Let  $\binom{n}{k} \equiv \frac{n!}{(n-k)!k!}$  where  $0! \equiv 1$  and  $(n+1)! \equiv (n+1)n!$  for all  $n \geq 0$ . Prove that whenever  $k \geq 1$  and  $k \leq n$ , then  $\binom{n+1}{k} = \binom{n}{k} + \binom{n}{k-1}$ . Are these numbers,  $\binom{n}{k}$  the same as those obtained in Pascal's triangle? Prove your assertion.
18. The binomial theorem states  $(a+b)^n = \sum_{k=0}^n \binom{n}{k} a^{n-k} b^k$ . Prove the binomial theorem by induction. **Hint:** You might try using the preceding problem.
19. Show that for  $p \in (0, 1)$ ,  $\sum_{k=0}^n \binom{n}{k} k p^k (1-p)^{n-k} = np$ .

20. Using the binomial theorem prove that for all  $n \in \mathbb{N}$ ,  $\left(1 + \frac{1}{n}\right)^n \leq \left(1 + \frac{1}{n+1}\right)^{n+1}$ .  
**Hint:** Show first that  $\binom{n}{k} = \frac{n \cdot (n-1) \cdots (n-k+1)}{k!}$ . By the binomial theorem,

$$\left(1 + \frac{1}{n}\right)^n = \sum_{k=0}^n \binom{n}{k} \left(\frac{1}{n}\right)^k = \sum_{k=0}^n \frac{\overbrace{n \cdot (n-1) \cdots (n-k+1)}^{k \text{ factors}}}{k! n^k}.$$

Now consider the term  $\frac{n \cdot (n-1) \cdots (n-k+1)}{k! n^k}$  and note that a similar term occurs in the binomial expansion for  $\left(1 + \frac{1}{n+1}\right)^{n+1}$  except that  $n$  is replaced with  $n+1$  wherever this occurs. Argue the term got bigger and then note that in the binomial expansion for  $\left(1 + \frac{1}{n+1}\right)^{n+1}$ , there are more terms.

21. Prove by induction that for all  $k \geq 4$ ,  $2^k \leq k!$
22. Use the Problems 21 and 20 to verify for all  $n \in \mathbb{N}$ ,  $\left(1 + \frac{1}{n}\right)^n \leq 3$ .
23. Prove by induction that  $1 + \sum_{i=1}^n i(i!) = (n+1)!$ .
24. I can jump off the top of the empire state building without suffering any ill effects. Here is the proof by induction. If I jump from a height of one inch, I am unharmed. Furthermore, if I am unharmed from jumping from a height of  $n$  inches, then jumping from a height of  $n+1$  inches will also not harm me. This is self evident and provides the induction step. Therefore, I can jump from a height of  $n$  inches for any  $n$ . What is the matter with this reasoning?

<sup>2</sup>Blaise Pascal lived in the 1600's and is responsible for the beginnings of the study of probability.



25. All horses are the same color. Here is the proof by induction. A single horse is the same color as himself. Now suppose the theorem that all horses are the same color is true for  $n$  horses and consider  $n + 1$  horses. Remove one of the horses and use the induction hypothesis to conclude the remaining  $n$  horses are all the same color. Put the horse which was removed back in and take out another horse. The remaining  $n$  horses are the same color by the induction hypothesis. Therefore, all  $n + 1$  horses are the same color as the  $n - 1$  horses which didn't get moved. This proves the theorem. Is there something wrong with this argument?

## 2.10 Divisibility And The Fundamental Theorem Of Arithmetic

It is not necessary to read this section in order to do calculus. However, it is good general knowledge and so is included. The following definition describes what is meant by a prime number and also what is meant by the word "divides".

**Definition 2.10.1** *The number,  $a$  divides the number,  $b$  if in Theorem 2.8.11,  $r = 0$ . That is there is zero remainder. The notation for this is  $a|b$ , read  $a$  divides  $b$  and  $a$  is called a factor of  $b$ . A prime number is one which has the property that the only numbers which divide it are itself and 1. The greatest common divisor of two positive integers,  $m, n$  is that number,  $p$  which has the property that  $p$  divides both  $m$  and  $n$  and also if  $q$  divides both  $m$  and  $n$ , then  $q$  divides  $p$ . Two integers are relatively prime if their greatest common divisor is one.*

**Theorem 2.10.2** *Let  $m, n$  be two positive integers and define*

$$S \equiv \{xm + yn \in \mathbb{N} : x, y \in \mathbb{Z}\}.$$

*Then the smallest number in  $S$  is the greatest common divisor, denoted by  $(m, n)$ .*

**Proof:** First note that both  $m$  and  $n$  are in  $S$  so it is a nonempty set of positive integers. By well ordering, there is a smallest element of  $S$ , called  $p = x_0m + y_0n$ . Either  $p$  divides  $m$  or it does not. If  $p$  does not divide  $m$ , then by Theorem 2.8.11,

$$m = pq + r$$

where  $0 < r < p$ . Thus  $m = (x_0m + y_0n)q + r$  and so, solving for  $r$ ,

$$r = m(1 - x_0) + (-y_0q)n \in S.$$

However, this is a contradiction because  $p$  was the smallest element of  $S$ . Thus  $p|m$ . Similarly  $p|n$ .

Now suppose  $q$  divides both  $m$  and  $n$ . Then  $m = qx$  and  $n = qy$  for integers,  $x$  and  $y$ . Therefore,

$$p = mx_0 + ny_0 = x_0qx + y_0qy = q(x_0x + y_0y)$$

showing  $q|p$ . Therefore,  $p = (m, n)$ .

**Theorem 2.10.3** *If  $p$  is a prime and  $p|ab$  then either  $p|a$  or  $p|b$ .*

**Proof:** Suppose  $p$  does not divide  $a$ . Then since the only factors of  $p$  are 1 and  $p$  it follows  $(p, a) = 1$  and therefore, there exists integers,  $x$  and  $y$  such that

$$1 = ax + yp.$$

Multiplying this equation by  $b$  yields

$$b = abx + ybp.$$

Since  $p|ab$ ,  $ab = pz$  for some integer  $z$ . Therefore,

$$b = abx + ybp = pzx + ybp = p(xz + yb)$$

and this shows  $p$  divides  $b$ .

**Theorem 2.10.4** (*Fundamental theorem of arithmetic*) Let  $a \in \mathbb{N} \setminus \{1\}$ . Then  $a = \prod_{i=1}^n p_i$  where  $p_i$  are all prime numbers. Furthermore, this prime factorization is unique except for the order of the factors.

**Proof:** If  $a$  equals a prime number, the prime factorization clearly exists. In particular the prime factorization exists for the prime number 2. Assume this theorem is true for all  $a \leq n-1$ . If  $n$  is a prime, then it has a prime factorization. On the other hand, if  $n$  is not a prime, then there exist two integers  $k$  and  $m$  such that  $n = km$  where each of  $k$  and  $m$  are less than  $n$ . Therefore, each of these is no larger than  $n-1$  and consequently, each has a prime factorization. Thus so does  $n$ . It remains to argue the prime factorization is unique except for order of the factors.

Suppose

$$\prod_{i=1}^n p_i = \prod_{j=1}^m q_j$$

where the  $p_i$  and  $q_j$  are all prime, there is no way to reorder the  $q_k$  such that  $m = n$  and  $p_i = q_i$  for all  $i$ , and  $n+m$  is the smallest positive integer such that this happens. Then by Theorem 2.10.3,  $p_1|q_j$  for some  $j$ . Since these are prime numbers this requires  $p_1 = q_1$ . Reordering if necessary it can be assumed that  $q_j = q_1$ . Then dividing both sides by  $p_1 = q_1$ ,

$$\prod_{i=1}^{n-1} p_{i+1} = \prod_{j=1}^{m-1} q_{j+1}.$$

Since  $n+m$  was as small as possible for the theorem to fail, it follows that  $n-1 = m-1$  and the prime numbers,  $q_2, \dots, q_m$  can be reordered in such a way that  $p_k = q_k$  for all  $k = 2, \dots, n$ . Hence  $p_i = q_i$  for all  $i$  because it was already argued that  $p_1 = q_1$ , and this results in a contradiction, proving the theorem.

The next theorem is a very nice high school theorem which characterizes all possible rational roots for polynomials having integer coefficients.

**Theorem 2.10.5** (*rational root theorem*) Let

$$a_n x^n + \dots + a_1 x + a_0 = 0$$

where each  $a_i$  is an integer and  $a_n \neq 0$ . Then if the equation has any rational solutions, these are of the form

$$\pm \frac{\text{factor of } a_0}{\text{factor of } a_n}.$$

**Proof:** Let  $\frac{p}{q}$  be a rational solution. Dividing  $p$  and  $q$  by  $(p, q)$  if necessary, the fraction may be reduced to lowest terms such that  $(p, q) = 1$ . Substituting into the equation,

$$a_n p^n + a_{n-1} p^{n-1} q + \dots + a_1 p q^{n-1} + a_0 q^n = 0.$$

Hence

$$a_n p^n = -(a_{n-1} p^{n-1} q + \cdots + a_0 q^n)$$

and  $q$  divides the right side of the equation and therefore,  $q$  must divide the left side also. However,  $(q, p^n) = 1$  and so by Theorem 2.10.3  $q|a_n$  because it does not divide  $p^n$  due to the fact that  $p^n$  and  $q$  have no prime factors in common.

Similarly,

$$a_0 q^n = -(a_n p^n + \cdots + a_1 p q^{n-1})$$

and so  $p|a_0 q^n$  but  $(p, q^n) = 1$ . By Theorem 2.10.3 again,  $p|a_0$  and this proves the theorem.

**Example 2.10.6** An irrational number is one which is not rational. Show  $\sqrt{2}$  is irrational if it exists.

$\sqrt{2}$  is the solution of the equation  $x^2 - 2 = 0$ . However, from Theorem 2.10.5, the only possible rational roots to this equation are  $\pm 2$  and  $\pm 1$  and none of these work. Therefore,  $\sqrt{2}$  must be irrational.

## 2.11 Exercises

- Using Theorem 2.10.5, show  $\sqrt[7]{6}$ ,  $\sqrt[3]{7}$ ,  $\sqrt[5]{5}$  are all irrational numbers. This means they are not rational.
- Using the fact that  $\sqrt{2}$  is irrational, (not rational) show that numbers of the form  $r\sqrt{2}$  where  $r \in \mathbb{Q}$  are dense in  $\mathbb{R}$ . Then verify these numbers are irrational.
- Euclid<sup>3</sup> showed there were infinitely many prime numbers using a very simple argument. He assumed there were only finitely many,  $\{p_1, \dots, p_n\}$  and then considered the number  $p_1 \cdots p_n + 1$  consisting of the product of all the primes plus 1. Then this number can't be prime because it is larger than every prime number. Therefore, some prime number,  $p_k$  from the above list must divide it. Now obtain a **terrible** contradiction.
- If  $a, b$  are integers,  $[a, b]$  will denote their least common multiple. This is the smallest number which has both  $a$  and  $b$  as factors. Show  $[a, b] = ab / (a, b)$ . **Hint:** Show  $[a, b]$  must divide  $ab$ . Here is how you might proceed. If not,  $ab = [a, b]q + r$  where  $0 < r < [a, b]$ . Then verify  $r$  is a common multiple of  $a$  and  $b$  contradicting that  $[a, b]$  is the **least** common multiple. Hence  $r = 0$ . Therefore,  $[a, b] = ab/q$  for some  $q$  an integer. Since  $[a, b]$  is a common multiple of  $a$  and  $b$ , argue that  $q$  must divide both  $a$  and  $b$ . Now what is the largest such  $q$ ? This would yield the smallest  $ab/q$ . You fill in the details.
- Show that if  $\{a, b, c\}$  are three positive integers, they have a greatest common divisor which may be written as  $ax + by + cz$  for some integers  $x, y, z$ .
- Let  $a_n = 2^{2^n} + 1$  for  $n = 1, 2, \dots$ . Show that if  $n \neq m$ , then  $a_n$  and  $a_m$  are relatively prime. Either  $a_n$  is prime or it is not. If it is not, then all the numbers dividing it other than 1 fail to divide  $a_m$  for all  $m < n$ . Explain why this shows there must be infinitely many primes. This argument about infinitely many primes is due to Polya. It gives more information than the argument of Euclid. The numbers,  $2^{(2^n)} + 1$  are

---

<sup>3</sup>He lived about 300 B.C.

prime numbers for several values of  $n$  but Euler<sup>4</sup> showed that when  $n = 5$ , the number is not prime<sup>5</sup>. When numbers of this form are prime, they are called Fermat<sup>6</sup> primes. At this time it is unknown whether there are infinitely many Fermat primes. For more information on these matters, you should see the book by Chahal, [4]. **Hint:** To verify  $a_n$  and  $a_m$  are relatively prime for  $m > n$ , suppose they are not and that for some number,  $p \neq 1$ ,  $a_n = pk_1$  while  $a_m = pk_2$ . Then letting  $m = n + r$ , explain why

$$\begin{aligned} pk_2 &= a_m = \left(2^{2^n}\right)^{2^r} + 1 = (pk_1 - 1)^{2^r} + 1 \\ &= p(\text{integer}) + 2. \end{aligned}$$

Consequently,  $p(\text{integer}) = 2$ . What does this say about  $p$ ? How does  $pk_1 = 2^{2^n} + 1$  yield a contradiction?

## 2.12 Systems Of Equations

Sometimes it is necessary to solve systems of equations. For example the problem could be to find  $x$  and  $y$  such that

$$x + y = 7 \text{ and } 2x - y = 8. \quad (2.3)$$

The set of ordered pairs,  $(x, y)$  which solve both equations is called the solution set. For example, you can see that  $(5, 2) = (x, y)$  is a solution to the above system. To solve this, note that the solution set does not change if any equation is replaced by a non zero multiple of itself. It also does not change if one equation is replaced by itself added to a multiple of the other equation. For example,  $x$  and  $y$  solve the above system if and only if  $x$  and  $y$  solve the system

$$x + y = 7, \overbrace{2x - y + (-2)(x + y) = 8 + (-2)(7)}^{-3y = -6}. \quad (2.4)$$

The second equation was replaced by  $-2$  times the first equation added to the second. Thus the solution is  $y = 2$ , from  $-3y = -6$  and now, knowing  $y = 2$ , it follows from the other equation that  $x + 2 = 7$  and so  $x = 5$ .

Why exactly does the replacement of one equation with a multiple of another added to it not change the solution set? The two equations of (2.3) are of the form

$$E_1 = f_1, E_2 = f_2 \quad (2.5)$$

where  $E_1$  and  $E_2$  are expressions involving the variables. The claim is that if  $a$  is a number, then (2.5) has the same solution set as

$$E_1 = f_1, E_2 + aE_1 = f_2 + af_1. \quad (2.6)$$

Why is this?

If  $(x, y)$  solves (2.5) then it solves the first equation in (2.6). Also, it satisfies  $aE_1 = af_1$  and so, since it also solves  $E_2 = f_2$  it must solve the second equation in (2.6). If  $(x, y)$

<sup>4</sup>Leonhard Euler, born in Switzerland, lived from 1707 to 1783. He was the most prolific mathematician ever to live. He made major contributions to number theory, analysis, algebra, mechanics, and differential equations. He and Lagrange invented the branch of mathematics known as calculus of variations. His collected papers take up more shelf space than a typical encyclopedia. His memory was prodigious and he could do unbelievable feats of computation in his head. He had 13 children.

<sup>5</sup>The number in this case is 4, 294, 967, 297.

<sup>6</sup>Fermat lived from 1601 to 1665. He is generally regarded as the founder of number theory. His most famous conjecture was that there is no solution to the equation  $x^n + y^n = z^n$  if  $n \geq 3$ . That is there is no analog to pythagorean triples with higher exponents than 2. This was finally proved in the 1990's by Andrew Wiles.

solves (2.6) then it solves the first equation of (2.5). Also  $aE_1 = af_1$  and it is given that the second equation of (2.6) is verified. Therefore,  $E_2 = f_2$  and it follows  $(x, y)$  is a solution of the second equation in (2.5). This shows the solutions to (2.5) and (2.6) are exactly the same which means they have the same solution set. Of course the same reasoning applies with no change if there are many more variables than two and many more equations than two. It is still the case that when one equation is replaced with a multiple of another one added to itself, the solution set of the whole system does not change.

The other thing which does not change the solution set of a system of equations consists of listing the equations in a different order. Here is another example.

**Example 2.12.1** Find the solutions to the system,

$$\begin{aligned}x + 3y + 6z &= 25 \\ 2x + 7y + 14z &= 58 \\ 2y + 5z &= 19\end{aligned}\tag{2.7}$$

To solve this system replace the second equation by  $(-2)$  times the first equation added to the second. This yields the system

$$\begin{aligned}x + 3y + 6z &= 25 \\ y + 2z &= 8 \\ 2y + 5z &= 19\end{aligned}\tag{2.8}$$

Now take  $(-2)$  times the second and add to the third. More precisely, replace the third equation with  $(-2)$  times the second added to the third. This yields the system

$$\begin{aligned}x + 3y + 6z &= 25 \\ y + 2z &= 8 \\ z &= 3\end{aligned}\tag{2.9}$$

At this point, you can tell what the solution is. This system has the same solution as the original system and in the above,  $z = 3$ . Then using this in the second equation, it follows  $y + 6 = 8$  and so  $y = 2$ . Now using this in the top equation yields  $x + 6 + 18 = 25$  and so  $x = 1$ .

This process is not really much different from what you have always done in solving a single equation. For example, suppose you wanted to solve  $2x + 5 = 3x - 6$ . You did the same thing to both sides of the equation thus preserving the solution set until you obtained an equation which was simple enough to give the answer. In this case, you would add  $-2x$  to both sides and then add 6 to both sides. This yields  $x = 11$ .

In (2.9) you could have continued as follows. Add  $(-2)$  times the bottom equation to the middle and then add  $(-6)$  times the bottom to the top. This yields

$$\begin{aligned}x + 3y &= 19 \\ y &= 6 \\ z &= 3\end{aligned}$$

Now add  $(-3)$  times the second to the top. This yields

$$\begin{aligned}x &= 1 \\ y &= 6, \\ z &= 3\end{aligned}$$

a system which has the same solution set as the original system.

It is foolish to write the variables every time you do these operations. It is easier to write the system (2.7) as the following “augmented matrix”

$$\begin{pmatrix} 1 & 3 & 6 & 25 \\ 2 & 7 & 14 & 58 \\ 0 & 2 & 5 & 19 \end{pmatrix}.$$

It has exactly the same information as the original system but here it is understood there is an  $x$  column,  $\begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}$ , a  $y$  column,  $\begin{pmatrix} 3 \\ 7 \\ 2 \end{pmatrix}$  and a  $z$  column,  $\begin{pmatrix} 6 \\ 14 \\ 5 \end{pmatrix}$ . The rows correspond to the equations in the system. Thus the top row in the augmented matrix corresponds to the equation,

$$x + 3y + 6z = 25.$$

Now when you replace an equation with a multiple of another equation added to itself, you are just taking a row of this augmented matrix and replacing it with a multiple of another row added to it. Thus the first step in solving (2.7) would be to take  $(-2)$  times the first row of the augmented matrix above and add it to the second row,

$$\begin{pmatrix} 1 & 3 & 6 & 25 \\ 0 & 1 & 2 & 8 \\ 0 & 2 & 5 & 19 \end{pmatrix}.$$

Note how this corresponds to (2.8). Next take  $(-2)$  times the second row and add to the third,

$$\begin{pmatrix} 1 & 3 & 6 & 25 \\ 0 & 1 & 2 & 8 \\ 0 & 0 & 1 & 3 \end{pmatrix}$$

which is the same as (2.9). You get the idea I hope. Write the system as an augmented matrix and follow the procedure of either switching rows, multiplying a row by a non zero number, or replacing a row by a multiple of another row added to it. Each of these operations leaves the solution set unchanged. These operations are called row operations.

**Example 2.12.2** Give the complete solution to the system of equations,  $5x + 10y - 7z = -2$ ,  $2x + 4y - 3z = -1$ , and  $3x + 6y + 5z = 9$ .

The augmented matrix for this system is

$$\begin{pmatrix} 2 & 4 & -3 & -1 \\ 5 & 10 & -7 & -2 \\ 3 & 6 & 5 & 9 \end{pmatrix}$$

Multiply the second row by 2, the first row by 5, and then take  $(-1)$  times the first row and add to the second. Then multiply the first row by  $1/5$ . This yields

$$\begin{pmatrix} 2 & 4 & -3 & -1 \\ 0 & 0 & 1 & 1 \\ 3 & 6 & 5 & 9 \end{pmatrix}$$

Now, combining some row operations, take  $(-3)$  times the first row and add this to 2 times the last row and replace the last row with this. This yields.

$$\begin{pmatrix} 2 & 4 & -3 & -1 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 21 \end{pmatrix}.$$

Putting in the variables, the last two rows say  $z = 1$  and  $z = 21$ . This is impossible so the last system of equations determined by the above augmented matrix has no solution. However, it has the same solution set as the first system of equations. This shows there is no solution to the three given equations. When this happens, the system is called inconsistent.

This should not be surprising that something like this can take place. It can even happen for one equation in one variable. Consider for example,  $x = x+1$ . There is clearly no solution to this.

**Example 2.12.3** Give the complete solution to the system of equations,  $3x - y - 5z = 9$ ,  $y - 10z = 0$ , and  $-2x + y = -6$ .

The augmented matrix of this system is

$$\left( \begin{array}{cccc} 3 & -1 & -5 & 9 \\ 0 & 1 & -10 & 0 \\ -2 & 1 & 0 & -6 \end{array} \right)$$

Replace the last row with 2 times the top row added to 3 times the bottom row. This gives

$$\left( \begin{array}{cccc} 3 & -1 & -5 & 9 \\ 0 & 1 & -10 & 0 \\ 0 & 1 & -10 & 0 \end{array} \right)$$

Next take  $-1$  times the middle row and add to the bottom.

$$\left( \begin{array}{cccc} 3 & -1 & -5 & 9 \\ 0 & 1 & -10 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right)$$

Take the middle row and add to the top and then divide the top row which results by 3.

$$\left( \begin{array}{cccc} 1 & 0 & -5 & 3 \\ 0 & 1 & -10 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right).$$

This says  $y = 10z$  and  $x = 3 + 5z$ . Apparently  $z$  can equal any number. Therefore, the solution set of this system is  $x = 3 + 5t$ ,  $y = 10t$ , and  $z = t$  where  $t$  is completely arbitrary. The system has an infinite set of solutions and this is a good description of the solutions. This is what it is all about, finding the solutions to the system.

The phenomenon of an infinite solution set occurs in equations having only one variable also. For example, consider the equation  $x = x$ . It doesn't matter what  $x$  equals.

**Definition 2.12.4** A system of linear equations is a list of equations,

$$\sum_{j=1}^n a_{ij}x_j = f_j, \quad i = 1, 2, 3, \dots, m$$

where  $a_{ij}$  are numbers,  $f_j$  is a number, and it is desired to find  $(x_1, \dots, x_n)$  solving each of the equations listed.

As illustrated above, such a system of linear equations may have a unique solution, no solution, or infinitely many solutions. It turns out these are the only three cases which can occur for linear systems. Furthermore, you do exactly the same things to solve any linear system. You write the augmented matrix and do row operations until you get a simpler system in which it is possible to see the solution. All is based on the observation that the row operations do not change the solution set. You can have more equations than variables, fewer equations than variables, etc. It doesn't matter. You always set up the augmented matrix and go to work on it. These things are all the same.

**Example 2.12.5** Give the complete solution to the system of equations,  $-41x + 15y = 168$ ,  $109x - 40y = -447$ ,  $-3x + y = 12$ , and  $2x + z = -1$ .

The augmented matrix is

$$\left( \begin{array}{cccc} -41 & 15 & 0 & 168 \\ 109 & -40 & 0 & -447 \\ -3 & 1 & 0 & 12 \\ 2 & 0 & 1 & -1 \end{array} \right).$$

To solve this multiply the top row by 109, the second row by 41, add the top row to the second row, and multiply the top row by  $1/109$ . This yields

$$\left( \begin{array}{cccc} -41 & 15 & 0 & 168 \\ 0 & -5 & 0 & -15 \\ -3 & 1 & 0 & 12 \\ 2 & 0 & 1 & -1 \end{array} \right).$$

Now take 2 times the third row and replace the fourth row by this added to 3 times the fourth row.

$$\left( \begin{array}{cccc} -41 & 15 & 0 & 168 \\ 0 & -5 & 0 & -15 \\ -3 & 1 & 0 & 12 \\ 0 & 2 & 3 & 21 \end{array} \right).$$

Take  $(-41)$  times the third row and replace the first row by this added to 3 times the first row. Then switch the third and the first rows.

$$\left( \begin{array}{cccc} 123 & -41 & 0 & -492 \\ 0 & -5 & 0 & -15 \\ 0 & 4 & 0 & 12 \\ 0 & 2 & 3 & 21 \end{array} \right).$$

Take  $-1/2$  times the third row and add to the bottom row. Then take 5 times the third row and add to four times the second. Finally take 41 times the third row and add to 4 times the top row. This yields

$$\left( \begin{array}{cccc} 492 & 0 & 0 & -1476 \\ 0 & 0 & 0 & 0 \\ 0 & 4 & 0 & 12 \\ 0 & 0 & 3 & 15 \end{array} \right)$$

It follows  $x = \frac{-1476}{492} = -3$ ,  $y = 3$  and  $z = 5$ .

You should practice solving systems of equations. Here are some exercises.

## 2.13 Exercises

1. Give the complete solution to the system of equations,  $3x - y + 4z = 6$ ,  $y + 8z = 0$ , and  $-2x + y = -4$ .
2. Give the complete solution to the system of equations,  $2x + z = 511$ ,  $x + 6z = 27$ , and  $y = 1$ .



3. Consider the system  $-5x + 2y - z = 0$  and  $-5x - 2y - z = 0$ . Both equations equal zero and so  $-5x + 2y - z = -5x - 2y - z$  which is equivalent to  $y = 0$ . Thus  $x$  and  $z$  can equal anything. But when  $x = 1$ ,  $z = -4$ , and  $y = 0$  are plugged in to the equations, it doesn't work. Why?
4. Give the complete solution to the system of equations,  $7x + 14y + 15z = 22$ ,  $2x + 4y + 3z = 5$ , and  $3x + 6y + 10z = 13$ .
5. Give the complete solution to the system of equations,  $-5x - 10y + 5z = 0$ ,  $2x + 4y - 4z = -2$ , and  $-4x - 8y + 13z = 8$ .
6. Give the complete solution to the system of equations,  $9x - 2y + 4z = -17$ ,  $13x - 3y + 6z = -25$ , and  $-2x - z = 3$ .
7. Give the complete solution to the system of equations,  $9x - 18y + 4z = -83$ ,  $-32x + 63y - 14z = 292$ , and  $-18x + 40y - 9z = 179$ .
8. Give the complete solution to the system of equations,  $65x + 84y + 16z = 546$ ,  $81x + 105y + 20z = 682$ , and  $84x + 110y + 21z = 713$ .
9. Give the complete solution to the system of equations,  $3x - y + 4z = -9$ ,  $y + 8z = 0$ , and  $-2x + y = 6$ .
10. Give the complete solution to the system of equations,  $8x + 2y + 3z = -3$ ,  $8x + 3y + 3z = -1$ , and  $4x + y + 3z = -9$ .
11. Give the complete solution to the system of equations,  $-7x - 14y - 10z = -17$ ,  $2x + 4y + 2z = 4$ , and  $2x + 4y - 7z = -6$ .
12. Give the complete solution to the system of equations,  $-8x + 2y + 5z = 18$ ,  $-8x + 3y + 5z = 13$ , and  $-4x + y + 5z = 19$ .
13. Give the complete solution to the system of equations,  $2x + 2y - 5z = 27$ ,  $2x + 3y - 5z = 31$ , and  $x + y - 5z = 21$ .
14. Give the complete solution to the system of equations,  $3x - y - 2z = 3$ ,  $y - 4z = 0$ , and  $-2x + y = -2$ .
15. Give the complete solution to the system of equations,  $3x - y - 2z = 6$ ,  $y - 4z = 0$ , and  $-2x + y = -4$ .
16. Four times the weight of Gaston is 150 pounds more than the weight of Ichabod. Four times the weight of Ichabod is 660 pounds less than seventeen times the weight of Gaston. Four times the weight of Gaston plus the weight of Siegfried equals 290 pounds. Brunhilde would balance all three of the others. Find the weights of the four girls.
17. Give the complete solution to the system of equations,  $-19x + 8y = -108$ ,  $-71x + 30y = -404$ ,  $-2x + y = -12$ ,  $4x + z = 14$ .
18. Give the complete solution to the system of equations,  $-9x + 15y = 66$ ,  $-11x + 18y = 79$ ,  $-x + y = 4$ , and  $z = 3$ .

## 2.14 Completeness of $\mathbb{R}$

By Theorem 2.8.9, between any two real numbers, points on the number line, there exists a rational number. This suggests there are a lot of rational numbers, but it is not clear from this Theorem whether the entire real line consists of only rational numbers. Some people might wish this were the case because then each real number could be described, not just as a point on a line but also algebraically, as the quotient of integers. Before 500 B.C., a group of mathematicians, led by Pythagoras believed in this, but they discovered their beliefs were false. It happened roughly like this. They knew they could construct the square root of two as the diagonal of a right triangle in which the two sides have unit length; thus they could regard  $\sqrt{2}$  as a number. Unfortunately, they were also able to show  $\sqrt{2}$  could not be written as the quotient of two integers. This discovery that the rational numbers could not even account for the results of geometric constructions was very upsetting to the Pythagoreans, especially when it became clear there were an endless supply of such “irrational” numbers.

This shows that if it is desired to consider all points on the number line, it is necessary to abandon the attempt to describe arbitrary real numbers in a purely algebraic manner using only the integers. Some might desire to throw out all the irrational numbers, and considering only the rational numbers, confine their attention to algebra, but this is not the approach to be followed here because it will effectively eliminate every major theorem of calculus. In this book real numbers will continue to be the points on the number line, a line which has no holes. This lack of holes is more precisely described in the following way.

**Definition 2.14.1** *A non empty set,  $S \subseteq \mathbb{R}$  is bounded above (below) if there exists  $x \in \mathbb{R}$  such that  $x \geq (\leq) s$  for all  $s \in S$ . If  $S$  is a nonempty set in  $\mathbb{R}$  which is bounded above, then a number,  $l$  which has the property that  $l$  is an upper bound and that every other upper bound is no smaller than  $l$  is called a least upper bound, l.u.b. ( $S$ ) or often  $\sup(S)$ . If  $S$  is a nonempty set bounded below, define the greatest lower bound, g.l.b. ( $S$ ) or  $\inf(S)$  similarly. Thus  $g$  is the g.l.b. ( $S$ ) means  $g$  is a lower bound for  $S$  and it is the largest of all lower bounds. If  $S$  is a nonempty subset of  $\mathbb{R}$  which is not bounded above, this information is expressed by saying  $\sup(S) = +\infty$  and if  $S$  is not bounded below,  $\inf(S) = -\infty$ .*

Every existence theorem in calculus depends on some form of the completeness axiom.

**Axiom 2.14.2** *(completeness) Every nonempty set of real numbers which is bounded above has a least upper bound and every nonempty set of real numbers which is bounded below has a greatest lower bound.*

It is this axiom which distinguishes Calculus from Algebra. A fundamental result about  $\sup$  and  $\inf$  is the following.

**Proposition 2.14.3** *Let  $S$  be a nonempty set and suppose  $\sup(S)$  exists. Then for every  $\delta > 0$ ,*

$$S \cap (\sup(S) - \delta, \sup(S)] \neq \emptyset.$$

*If  $\inf(S)$  exists, then for every  $\delta > 0$ ,*

$$S \cap [\inf(S), \inf(S) + \delta) \neq \emptyset.$$

**Proof:** Consider the first claim. If the indicated set equals  $\emptyset$ , then  $\sup(S) - \delta$  is an upper bound for  $S$  which is smaller than  $\sup(S)$ , contrary to the definition of  $\sup(S)$  as the least upper bound. In the second claim, if the indicated set equals  $\emptyset$ , then  $\inf(S) + \delta$  would be a lower bound which is larger than  $\inf(S)$  contrary to the definition of  $\inf(S)$ .

## 2.15 Review Exercises

1. Let  $S = [2, 5]$ . Find  $\sup S$ . Now let  $S = [2, 5)$ . Find  $\sup S$ . Is  $\sup S$  always a number in  $S$ ? Give conditions under which  $\sup S \in S$  and then give conditions under which  $\inf S \in S$ .
2. Show that if  $S \neq \emptyset$  and is bounded above (below) then  $\sup S$  ( $\inf S$ ) is unique. That is, there is only one least upper bound and only one greatest lower bound. If  $S = \emptyset$  can you conclude that 7 is an upper bound? Can you conclude 7 is a lower bound? What about 13.5? What about any other number?
3. Let  $S$  be a set which is bounded above and let  $-S$  denote the set  $\{-x : x \in S\}$ . How are  $\inf(-S)$  and  $\sup(S)$  related? **Hint:** Draw some pictures on a number line. What about  $\sup(-S)$  and  $\inf S$  where  $S$  is a set which is bounded below?

4. Solve the following equations which involve absolute values.

$$(a) |x + 1| = |2x + 3|$$

$$(b) |x + 1| - |x + 4| = 6$$

5. Solve the following inequalities which involve absolute values.

$$(a) |2x - 6| < 4$$

$$(b) |x - 2| < |2x + 2|$$

6. Which of the field axioms is being abused in the following argument that  $0 = 2$ ? Let  $x = y = 1$ . Then

$$0 = x^2 - y^2 = (x - y)(x + y)$$

and so

$$0 = (x - y)(x + y).$$

Now divide both sides by  $x - y$  to obtain

$$0 = x + y = 1 + 1 = 2.$$

7. Give conditions under which equality holds in the triangle inequality.
8. Let  $k \leq n$  where  $k$  and  $n$  are natural numbers.  $P(n, k)$ , permutations of  $n$  things taken  $k$  at a time, is defined to be the number of different ways to form an ordered list of  $k$  of the numbers,  $\{1, 2, \dots, n\}$ . Show

$$P(n, k) = n \cdot (n - 1) \cdots (n - k + 1) = \frac{n!}{(n - k)!}.$$

9. Using the preceding problem, show the number of ways of selecting a set of  $k$  things from a set of  $n$  things is  $\binom{n}{k}$ .
10. Prove the binomial theorem from Problem 9. **Hint:** When you take  $(x + y)^n$ , note that the result will be a sum of terms of the form,  $a_k x^{n-k} y^k$  and you need to determine what  $a_k$  should be. Imagine writing  $(x + y)^n = (x + y)(x + y) \cdots (x + y)$  where there are  $n$  factors in the product. Now consider what happens when you multiply. Each factor contributes either an  $x$  or a  $y$  to a typical term.
11. Prove by induction that  $n < 2^n$  for all natural numbers,  $n \geq 1$ .

12. Prove by the binomial theorem and Problem 9 that the number of subsets of a given finite set containing  $n$  elements is  $2^n$ .
13. Let  $n$  be a natural number and let  $k_1 + k_2 + \cdots + k_r = n$  where  $k_i$  is a non negative integer. The symbol

$$\binom{n}{k_1 k_2 \cdots k_r}$$

denotes the number of ways of selecting  $r$  subsets of  $\{1, \dots, n\}$  which contain  $k_1, k_2, \dots, k_r$  elements in them. Find a formula for this number.

14. Is it ever the case that  $(a+b)^n = a^n + b^n$  for  $a$  and  $b$  positive real numbers?
15. Is it ever the case that  $\sqrt{a^2 + b^2} = a + b$  for  $a$  and  $b$  positive real numbers?
16. Is it ever the case that  $\frac{1}{x+y} = \frac{1}{x} + \frac{1}{y}$  for  $x$  and  $y$  positive real numbers?
17. Derive a formula for the multinomial expansion,  $(\sum_{k=1}^p a_k)^n$  which is analogous to the binomial expansion. **Hint:** See Problem 10.
18. Suppose  $a > 0$  and that  $x$  is a real number which satisfies the quadratic equation,

$$ax^2 + bx + c = 0.$$

Find a formula for  $x$  in terms of  $a$  and  $b$  and square roots of expressions involving these numbers. **Hint:** First divide by  $a$  to get

$$x^2 + \frac{b}{a}x + \frac{c}{a} = 0.$$

Then add and subtract the quantity  $b^2/4a^2$ . Verify that

$$x^2 + \frac{b}{a}x + \frac{b^2}{4a^2} = \left(x + \frac{b}{2a}\right)^2.$$

Now solve the result for  $x$ . The process by which this was accomplished in adding in the term  $b^2/4a^2$  is referred to as completing the square. You should obtain the quadratic formula<sup>7</sup>,

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}.$$

The expression  $b^2 - 4ac$  is called the discriminant. When it is positive there are two different real roots. When it is zero, there is exactly one real root and when it equals a negative number there are no real roots.

19. Suppose  $f(x) = 3x^2 + 7x - 17$ . Find the value of  $x$  at which  $f(x)$  is smallest by completing the square. Also determine  $f(\mathbb{R})$  and sketch the graph of  $f$ . **Hint:**

$$\begin{aligned} f(x) &= 3\left(x^2 + \frac{7}{3}x - \frac{17}{3}\right) = 3\left(x^2 + \frac{7}{3}x + \frac{49}{36} - \frac{49}{36} - \frac{17}{3}\right) \\ &= 3\left(\left(x + \frac{7}{6}\right)^2 - \frac{49}{36} - \frac{17}{3}\right). \end{aligned}$$

---

<sup>7</sup>The ancient Babylonians knew how to solve these quadratic equations sometime before 1700 B.C. It seems they used pretty much the same process outlined in this exercise.

20. Suppose  $f(x) = -5x^2 + 8x - 7$ . Find  $f(\mathbb{R})$ . In particular, find the largest value of  $f(x)$  and the value of  $x$  at which it occurs. Can you conjecture and prove a result about  $y = ax^2 + bx + c$  in terms of the sign of  $a$  based on these last two problems?
21. Show that if it is assumed  $\mathbb{R}$  is complete, then the Archimedian property can be proved.  
**Hint:** Suppose completeness and let  $a > 0$ . If there exists  $x \in \mathbb{R}$  such that  $na \leq x$  for all  $n \in \mathbb{N}$ , then  $x/a$  is an upper bound for  $\mathbb{N}$ . Let  $l$  be the least upper bound and argue there exists  $n \in \mathbb{N} \cap [l - 1/4, l]$ . Now what about  $n + 1$ ?

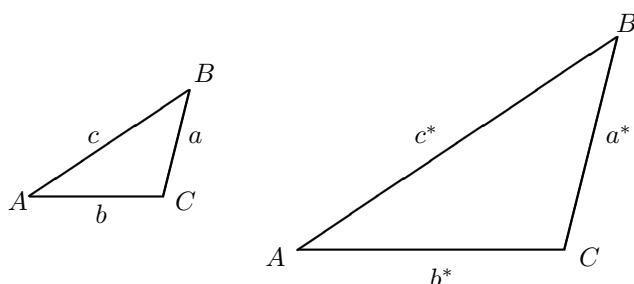


# Basic Geometry And Trigonometry

This section is a review some basic geometry which is especially useful in the study of calculus. The purpose here is not to give a complete treatment of plane geometry, just a suitable introduction. To do this right, you should consult the books of Euclid written about 300 B.C. [6]

## 3.1 Similar Triangles And Pythagorean Theorem

**Definition 3.1.1** *Two triangles are similar if they have the same angles. For example, in the following picture, the two triangles are similar because the angles are the same.*



The fundamental axiom for similar triangles is the following.

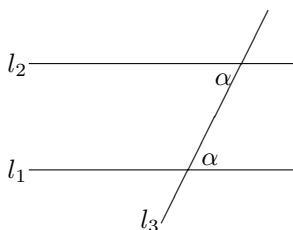
**Axiom 3.1.2** *If two triangles are similar then the ratios of corresponding parts are the same.*

For example in the above picture, this says that

$$\frac{a}{b} = \frac{a^*}{b^*}$$

**Definition 3.1.3** *Two lines in the plane are said to be parallel if no matter how far they are extended, they never intersect.*

**Definition 3.1.4** *If two lines  $l_1$  and  $l_2$  are parallel and if they are intersected by a line,  $l_3$ , the alternate interior angles are shown in the following picture labeled as  $\alpha$ .*



As suggested by the above picture, the following axiom will be used.

**Axiom 3.1.5** *If  $l_1$  and  $l_2$  are parallel lines intersected by  $l_3$ , then alternate interior angles are equal.*

**Definition 3.1.6** *An angle is a right angle if when either side is extended, the new angle formed by the extension equals the original angle.*

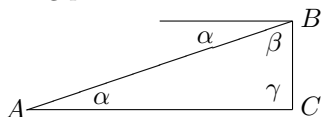
**Axiom 3.1.7** *Suppose  $l_1$  and  $l_2$  both intersect a third line,  $l_3$  in a right angle. Then  $l_1$  and  $l_2$  are parallel.*

**Definition 3.1.8** *A right triangle is one in which one of the angles is a right angle.*

**Axiom 3.1.9** *Given a straight line and a point, there exists a straight line which contains the point and intersects the given line in two right angles. This line is called perpendicular to the given line.*

**Theorem 3.1.10** *Let  $\alpha, \beta$ , and  $\gamma$  be the angles of a right triangle with  $\gamma$  the right angle. Then if the angles,  $\alpha$  and  $\beta$  are placed next to each other, the resulting angle is a right angle.*

**Proof:** Consider the following picture.



In the picture the top horizontal line is obtained from Axiom 3.1.9. It is a line perpendicular to the line determined by the line segment joining  $B$  and  $C$  which passes through the point,  $B$ . Then as shown in the picture, the angle formed by placing  $\alpha$  and  $\beta$  together is a right angle as claimed.

**Definition 3.1.11** *When an angle  $\alpha$  is placed next to an angle  $\beta$  as shown above, then the resulting angle is denoted by  $\alpha + \beta$ . A right angle is said to have  $90^\circ$  or to be a  $90^\circ$  angle.*

With this definition, Theorem 3.1.10 says the sum of the two non  $90^\circ$  angles in a right triangle is  $90^\circ$ .

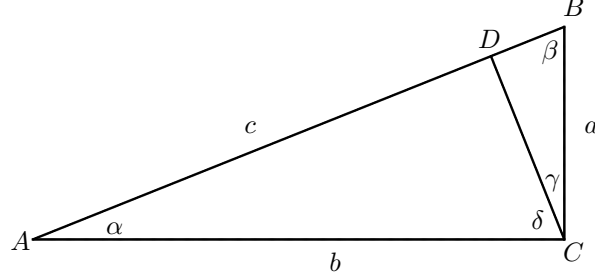
In a right triangle the long side is called the hypotenuse. The similar triangles axiom can be used to prove the Pythagorean theorem.

**Theorem 3.1.12 (Pythagoras)** *In a right triangle the square of the length of the hypotenuse equals the sum of the squares of the lengths of the other two sides.*

**Proof:** Consider the following picture in which the large triangle is a right triangle and  $D$  is the point where the line through  $C$  perpendicular to the line from  $A$  to  $B$  intersects the line from  $A$  to  $B$ . Then  $c$  is defined to be the length of the line from  $A$  to  $B$ ,  $a$  is the



length of the line from  $B$  to  $C$ , and  $b$  is the length of the line from  $A$  to  $C$ . Denote by  $\overline{DB}$  the length of the line from  $D$  to  $B$ .



Then from Theorem 3.1.10,  $\delta + \gamma = 90^\circ$  and  $\beta + \gamma = 90^\circ$ . Therefore,  $\delta = \beta$ . Also from this same theorem,  $\alpha + \delta = 90^\circ$  and so  $\alpha = \gamma$ . Therefore, the three triangles shown in the picture are all similar. By Axiom 3.1.2,

$$\frac{c}{a} = \frac{a}{\overline{DB}}, \text{ and } \frac{c}{b} = \frac{b}{c - \overline{DB}}.$$

Therefore,  $c\overline{DB} = a^2$  and

$$c(c - \overline{DB}) = b^2$$

so

$$\begin{aligned} c^2 &= c\overline{DB} + b^2 \\ &= a^2 + b^2. \end{aligned}$$

This proves the Pythagorean theorem. <sup>1</sup>

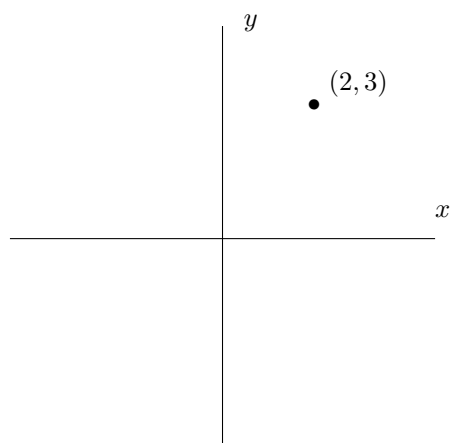
This theorem implies there should exist some such number which deserves to be called  $\sqrt{a^2 + b^2}$  as mentioned earlier in the discussion on completeness of  $\mathbb{R}$ .

## 3.2 Cartesian Coordinates And Straight Lines

Recall the notion of the Cartesian coordinate system. It involved an  $x$  axis, a  $y$  axis, two lines which intersect each other at right angles and one identifies a point by specifying a pair of numbers. For example, the number  $(2, 3)$  involves going 2 units to the right on the  $x$  axis and then 3 units directly up on a line perpendicular to the  $x$  axis. For example, consider the following picture.

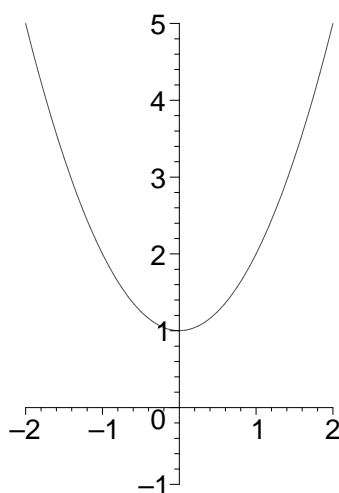
---

<sup>1</sup>This theorem is due to Pythagoras who lived about 572-497 B.C. This was during the Babylonian captivity of the Jews. Thus Pythagoras was probably a contemporary of the prophet Daniel, sometime before Ezra and Nehemiah. Alexander the great would not come along for more than 100 years. There was, however, an even earlier Greek mathematician named Thales, 624-547 B.C. who also did fundamental work in geometry. Greek geometry was organized and published by Euclid about 300 B.C.

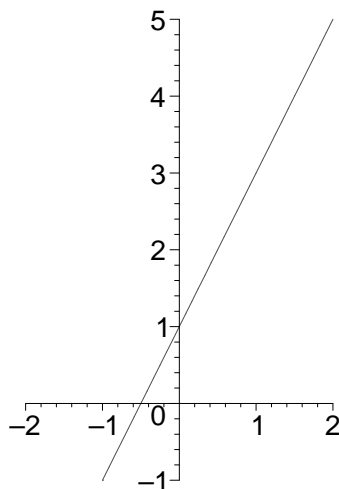


Because of the simple correspondence between points in the plane and the coordinates of a point in the plane, it is often the case that people are a little sloppy in referring to these things. Thus, it is common to see  $(x, y)$  referred to as a point in the plane. I will often indulge in this sloppiness.

The reader has likely encountered the notion of graphing relations of the form  $y = 2x + 3$  or  $y = x^2 + 5$ . Recall that you first found lots of ordered pairs which satisfied the relation. For example  $(0, 3)$ ,  $(1, 5)$ , and  $(-1, 1)$  all satisfy the first relation which describes a straight line. Here are some simple examples which you should see that you understand. First here is the graph of  $y = x^2 + 1$ .



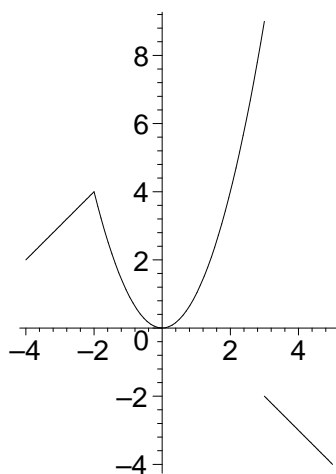
Now here is the graph of the relation  $y = 2x + 1$  which is a straight line.



Sometimes a relation is defined using different formulas depending on the location of one of the variables. For example, consider

$$y = \begin{cases} 6 + x & \text{if } x \leq -2 \\ x^2 & \text{if } -2 < x < 3 \\ 1 - x & \text{if } x \geq 3 \end{cases}$$

Then the graph of this relation is sketched below.



A very important type of relation is one of the form  $y - y_0 = m(x - x_0)$ , where  $m, x_0$ , and  $y_0$  are numbers. The reason this is important is that if there are two points,  $(x_1, y_1)$ , and  $(x_2, y_2)$  which satisfy this relation, then

$$\begin{aligned} \frac{y_1 - y_2}{x_1 - x_2} &= \frac{(y_1 - y_0) - (y_2 - y_0)}{x_1 - x_2} = \frac{m(x_1 - x_0) - m(x_2 - x_0)}{x_1 - x_2} \\ &= \frac{m(x_1 - x_2)}{x_1 - x_2} = m. \end{aligned}$$

Remember the slope of the line segment through two points is always the difference in the  $y$  values divided by the difference in the  $x$  values, taken in the same order. Sometimes this

is referred to as the rise divided by the run. This shows that there is a constant slope,  $m$ , the slope of the line, between any pair of points satisfying this relation. Such a relation is called a straight line. Also, the point  $(x_0, y_0)$  satisfies the relation. This is more often called the equation of the straight line.

**Example 3.2.1** Find the relation for a straight line which contains the point  $(1, 2)$  and has constant slope equal to 3.

From the above discussion,  $(y - 2) = 3(x - 1)$ .

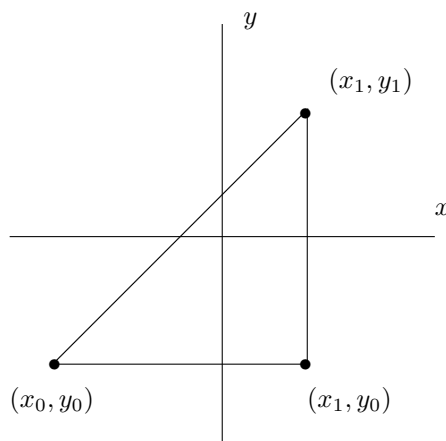
### 3.3 Exercises

1. Sketch the graph of  $y = x^3 + 1$ .
2. Sketch the graph of  $y = x^2 - 2x + 1$ .
3. Sketch the graph of  $y = \frac{x}{x^2 + 1}$ .
4. Sketch the graph of  $\frac{1}{1 + x^2}$ .
5. Sketch the graph of the straight line which goes through the points  $(1, 0)$  and  $(2, 3)$  and find the relation which describes this line.
6. Suppose  $a, b \neq 0$ . Find the equation of the line which goes through the points  $(0, a)$ , and  $(b, 0)$ .
7. Two lines are parallel if they have the same slope. Find the equation of the line through the point  $(2, 3)$  which is parallel to the line whose equation is  $2x + 3y = 8$ .
8. Sketch the graph of the relation defined as

$$y = \begin{cases} 1 & \text{if } x \leq -2 \\ 1 - x & \text{if } -2 < x < 3 \\ 1 + x & \text{if } x \geq 3 \end{cases}$$

### 3.4 Distance Formula And Trigonometric Functions

As just explained, points in the plane may be identified by giving a pair of numbers. Suppose there are two points in the plane and it is desired to find the distance between them. There are actually many ways used to measure this distance but the best way, and the only way used in this book is determined by the Pythagorean theorem. Consider the following picture.



In this picture, the distance between the points denoted by  $(x_0, y_0)$  and  $(x_1, y_1)$  should be the square root of the sum of the squares of the lengths of the two sides. The length of the side on the bottom is  $|x_0 - x_1|$  while the length of the side on the right is  $|y_0 - y_1|$ . Therefore, by the Pythagorean theorem the distance between the two indicated points is  $\sqrt{(x_0 - x_1)^2 + (y_0 - y_1)^2}$ . Note you could write

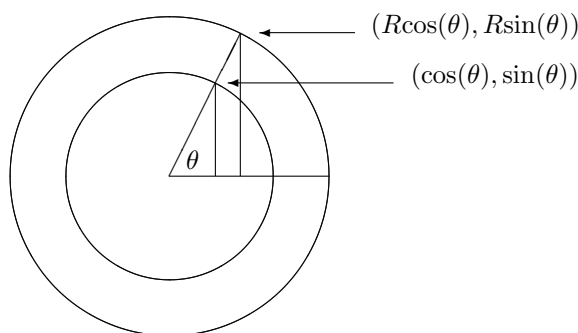
$$\sqrt{(x_1 - x_0)^2 + (y_1 - y_0)^2}$$

or even

$$\sqrt{(x_0 - x_1)^2 + (y_1 - y_0)^2}$$

and it would make no difference in the resulting number. The distance between the two points is written as  $|(x_0, y_0) - (x_1, y_1)|$  or sometimes when  $P_0$  is the point determined by  $(x_0, y_0)$  and  $P_1$  is the point determined by  $(x_1, y_1)$ , as  $d(P_0, P_1)$  or  $|P_0P_1|$ .

The trigonometric functions  $\cos$  and  $\sin$  are defined next. Consider the following picture in which the small circle has radius 1, the large circle has radius  $R$ , and the right side of each of the two triangles is perpendicular to the bottom side which lies on the  $x$  axis.



By Theorem 3.1.10 on Page 48 the two triangles have the same angles and so they are similar. Now define by  $(\cos \theta, \sin \theta)$  the coordinates of the top vertex of the smaller triangle. Therefore, it follows the coordinates of the top vertex of the larger triangle are as shown. This shows the following definition is well defined.

**Definition 3.4.1** For  $\theta$  an angle, define  $\cos \theta$  and  $\sin \theta$  as follows. Place the vertex of the angle (The vertex is the point.) at the point whose coordinates are  $(0,0)$  in such a way that one side of the angle lies on the positive  $x$  axis and the other side extends upward. Extend this other side until it intersects a circle of radius  $R$ . Then the point of intersection, is given as  $(R \cos \theta, R \sin \theta)$ . In particular, this specifies  $\cos \theta$  and  $\sin \theta$  by simply letting  $R = 1$ .

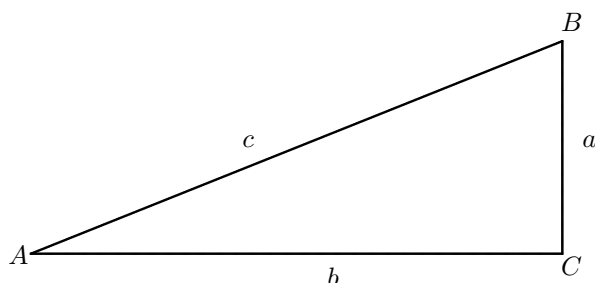
**Proposition 3.4.2** For any angle,  $\theta$ ,  $\cos^2 \theta + \sin^2 \theta = 1$ .

**Proof:** This follows immediately from the above definition and the distance formula. Since  $(\cos \theta, \sin \theta)$  is a point on the circle which has radius 1, the distance of this point to  $(0,0)$  equals 1. Thus the above identity holds.

The other trigonometric functions are defined as follows.

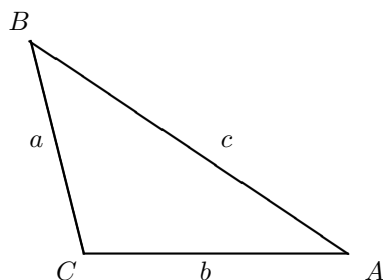
$$\tan \theta \equiv \frac{\sin \theta}{\cos \theta}, \cot \theta \equiv \frac{\cos \theta}{\sin \theta}, \sec \theta \equiv \frac{1}{\cos \theta}, \csc \theta \equiv \frac{1}{\sin \theta}. \quad (3.1)$$

It is also important to understand these functions in terms of a right triangle. Consider the following picture of a right triangle.



You should verify  $\sin A \equiv a/c$ ,  $\cos A \equiv b/c$ ,  $\tan A \equiv a/b$ ,  $\sec A \equiv c/b$ , and  $\csc A \equiv c/a$ .

Having defined the  $\cos$  and  $\sin$  there is a very important generalization of the Pythagorean theorem known as the law of cosines. Consider the following picture of a triangle in which  $a, b$  and  $c$  are the lengths of the sides and  $A, B$ , and  $C$  denote the angles indicated.



The law of cosines is the following.

**Theorem 3.4.3** Let  $ABC$  be a triangle as shown above. Then

$$c^2 = a^2 + b^2 - 2ab \cos C$$

**Proof:** Situate the triangle so the vertex of the angle,  $C$ , is on the point whose coordinates are  $(0,0)$  and so the side opposite the vertex,  $B$  is on the positive  $x$  axis as shown in the above picture. Then from the definition of the  $\cos C$ , the coordinates of the vertex,

$B$  are  $(a \cos C, a \sin C)$  while it is clear the coordinates of  $A$  are  $(b, 0)$ . Therefore, from the distance formula, and Proposition 3.4.2,

$$\begin{aligned} c^2 &= (a \cos C - b)^2 + a^2 \sin^2 C \\ &= a^2 \cos^2 C - 2ab \cos C + b^2 + a^2 \sin^2 C \\ &= a^2 + b^2 - 2ab \cos C \end{aligned}$$

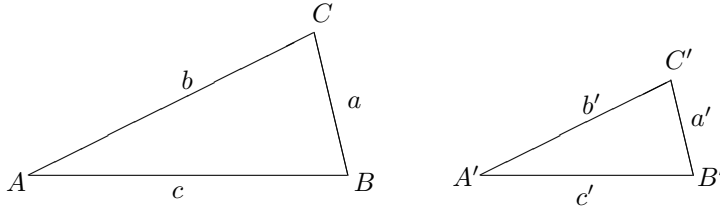
as claimed.

**Corollary 3.4.4** *Let  $ABC$  be any triangle as shown above. Then the length of any side is no longer than the sum of the lengths of the other two sides.*

**Proof:** This follows immediately from the law of cosines. From Proposition 3.4.2,  $|\cos \theta| \leq 1$  and so  $c^2 = a^2 + b^2 - 2ab \cos C \leq a^2 + b^2 + 2ab = (a + b)^2$ . This proves the corollary.

**Corollary 3.4.5** *Suppose  $T$  and  $T'$  are two triangles such that one angle is the same in the two triangles and in each triangle, the sides forming that angle are equal. Then the corresponding sides are proportional.*

**Proof:** Let  $T = ABC$  with the two equal sides being  $AC$  and  $AB$ . Let  $T'$  be labeled in the same way but with primes on the letters. Thus the angle at  $A$  is equal to the angle at  $A'$ . The following picture is descriptive of the situation.



Denote by  $a, a', b, b', c$  and  $c'$  the sides indicated in the picture. Then by the law of cosines,

$$\begin{aligned} a^2 &= b^2 + c^2 - 2bc \cos A \\ &= 2b^2 - 2b^2 \cos A \end{aligned}$$

and so  $a/b = \sqrt{2(1 - \cos A)}$ . Similar reasoning shows  $a'/b' = \sqrt{2(1 - \cos A)}$  and so

$$a/b = a'/b'.$$

Similarly,  $a/c = a'/c'$ . By assumption  $c/b = 1 = c'/b'$ .

Such triangles in which two sides are equal are called isosceles.

### 3.5 The Circular Arc Subtended By An Angle

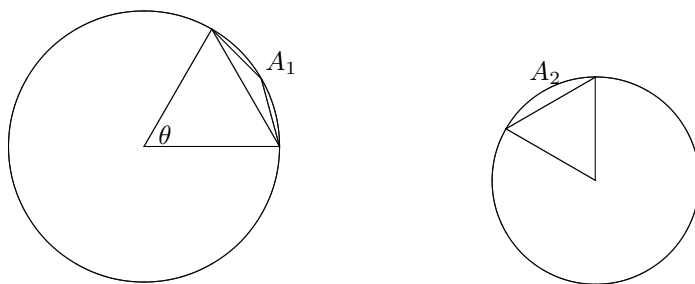
How can angles be measured? This will be done by considering arcs on a circle. To see how this will be done, let  $\theta$  denote an angle and place the vertex of this angle at the center of

the circle. Next, extend its two sides till they intersect the circle. Note the angle could be opening in any of infinitely many different directions. Thus this procedure could yield any of infinitely many different circular arcs. Each of these arcs is said to subtend the angle. In fact each of these arcs has the same length. When this has been shown, it will be easy to measure angles. Angles will be measured in terms of lengths of arcs subtended by the angle. Of course it is also necessary to define what is meant by the length of a circular arc in order to do any of this. First I will describe an intuitive way of thinking about this and then give a rigorous definition and proof. If the intuitive way of thinking about this satisfies you, no harm will be done by skipping the more technical discussion which follows.

Take an angle and place its vertex (the point) at the center of a circle of radius  $r$ . Then, extending the sides of the angle if necessary till they intersect the circle, this determines an arc on the circle. If  $r$  were changed to  $R$ , this really amounts to a change of units of length. Think, for example, of keeping the numbers the same but changing centimeters to meters in order to produce an enlarged version of the same picture. Thus the picture looks exactly the same, only larger. It is reasonable to suppose, based on this reasoning that the way to measure the angle is to take the length of the arc subtended in whatever units being used and divide this length by the radius measured in the same units, thus obtaining a number which is independent of the units of length used just as the angle itself is independent of units of length. After all, it is the same angle regardless of how far its sides are extended. This is in fact how to define the radian measure of an angle and the definition is well defined. Thus, in particular, the ratio between the circumference (length) of a circle and its radius is a constant which is independent of the radius of the circle<sup>2</sup>. Since the time of Euler in the 1700's, this constant has been denoted by  $2\pi$ . In summary, if  $\theta$  is the radian measure of an angle, the length of the arc subtended by the angle on a circle of radius  $r$  is  $r\theta$ .

This is a little sloppy right now because no precise definition of the length of an arc of a circle has been given. For now, imagine taking a string, placing one end of it on one end of the circular arc and then wrapping the string till you reach the other end of the arc. Stretching this string out and measuring it would then give you the length of the arc. Such intuitive discussions involving string may or may not be enough to convey understanding. If you need to see more discussion, read on. Otherwise, skip to the next section.

To give a precise description of what is meant by the length of an arc, consider the following picture.



In this picture, there are two circles, a big one having radius,  $R$  and a little one having radius  $r$ . The angle,  $\theta$  is situated in two different ways subtending the arcs  $A_1$  and  $A_2$  as shown.

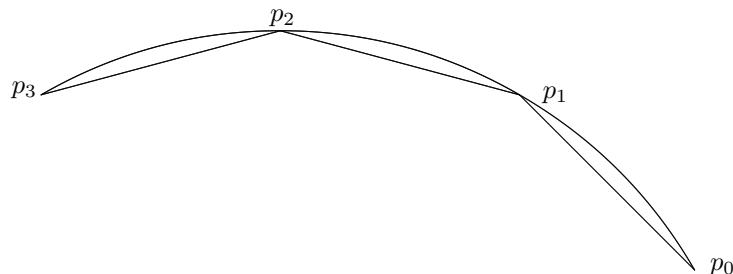
Letting  $A$  be an arc of a circle, like those shown in the above picture, A subset of

---

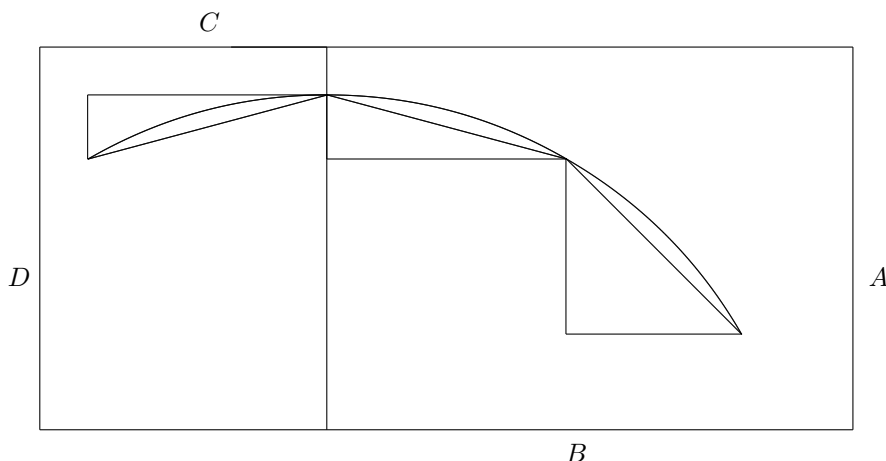
<sup>2</sup>In 2 Chronicles 4:2 the "molten sea" used for "washing" by the priests and found in Solomon's temple is described. It sat on 12 oxen, was round, 5 cubits high, 10 across and 30 around. Thus the Bible, taken literally, gives the value of  $\pi$  as 3. This is not too far off. Later, methods will be given which allow one to calculate  $\pi$  more precisely. A better value is 3.1415926535 and presently this number is known to thousands of decimal places. It was proved by Lindeman in the 1880's that  $\pi$  is transcendental which is the worst sort of irrational number.



$A, \{p_0, \dots, p_n\}$  is a partition of  $A$  if  $p_0$  is one endpoint,  $p_n$  is the other end point, and the points are encountered in the indicated order as one moves in the counter clockwise direction along the arc. To illustrate, see the following picture.

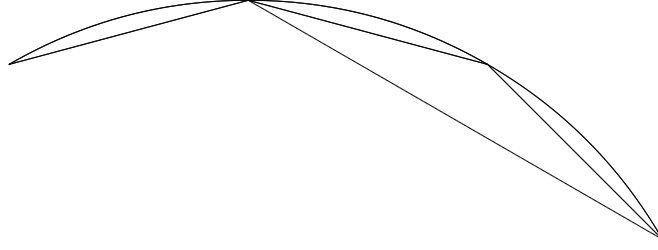


Also, denote by  $\mathcal{P}(A)$  the set of all such partitions. For  $P = \{p_0, \dots, p_n\}$ , denote by  $|p_i - p_{i-1}|$  the distance between  $p_i$  and  $p_{i-1}$ . Then for  $P \in \mathcal{P}(A)$ , define  $|P| \equiv \sum_{i=1}^n |p_i - p_{i-1}|$ . Thus  $|P|$  consists of the sum of the lengths of the little lines joining successive points of  $P$  and appears to be an approximation to the length of the circular arc,  $A$ . By Corollary 3.4.4 the length of any of the straight line segments joining successive points in a partition is smaller than the sum of the two sides of a right triangle having the given straight line segment as its hypotenuse. For example, see the following picture.

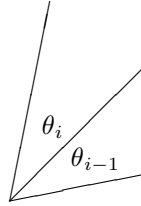


The sum of the lengths of the straight line segments in the part of the picture found in the right rectangle above is less than  $A + B$  and the sum of the lengths of the straight line segments in the part of the picture found in the left rectangle above is less than  $C + D$  and this would be so for any partition. Therefore, for any  $P \in \mathcal{P}(A)$ ,  $|P| \leq M$  where  $M$  is the perimeter of a rectangle containing the arc,  $A$ . To be a little sloppy, simply pick  $M$  to be the perimeter of a rectangle containing the whole circle of which  $A$  is a part. The only purpose for doing this is to obtain the existence of an upper bound. Therefore,  $\{|P| : P \in \mathcal{P}(A)\}$  is a set of numbers which is bounded above by  $M$  and by completeness of  $\mathbb{R}$  it is possible to define the length of  $A$ ,  $l(A)$ , by  $l(A) \equiv \sup \{|P| : P \in \mathcal{P}(A)\}$ .

A fundamental observation following from Corollary 3.4.4 is that if  $P, Q \in \mathcal{P}(A)$  and  $P \subseteq Q$ , then  $|P| \leq |Q|$ . To see this, add in one point at a time to  $P$ . This effect of adding in one point is illustrated in the following picture.



Also, letting  $\{p_0, \dots, p_n\}$  be a partition of  $A$ , specify angles,  $\theta_i$  as follows. The angle  $\theta_i$  is formed by the two lines, one from the center of the circle to  $p_i$  and the other line from the center of the circle to  $p_{i-1}$ . Furthermore, a specification of these angles yields the partition of  $A$  in the following way. Place the vertex of  $\theta_1$  on the center of the circle, letting one side lie on the line from the center of the circle to  $p_0$  and the other side extended resulting in a point further along the arc in the counter clockwise direction. When the angles,  $\theta_1, \dots, \theta_{i-1}$  have produced points,  $p_0, \dots, p_{i-1}$  on the arc, place the vertex of  $\theta_i$  on the center of the circle and let one side of  $\theta_i$  coincide with the side of the angle  $\theta_{i-1}$  which is most counter clockwise, the other side of  $\theta_i$  when extended, resulting in a point further along the arc,  $A$  in the counterclockwise direction as shown below.



Now let  $\varepsilon > 0$  be given and pick  $P_1 \in \mathcal{P}(A_1)$  such that  $|P_1| + \varepsilon > l(A_1)$ . Then determining the angles as just described, use these angles to produce a corresponding partition of  $A_2$ ,  $P_2$ . If  $|P_2| + \varepsilon > l(A_2)$ , then stop. Otherwise, pick  $Q \in \mathcal{P}(A_2)$  such that  $|Q| + \varepsilon > l(A_2)$  and let  $P'_2 = P_2 \cup Q$ . Then use the angles determined by  $P'_2$  to obtain  $P'_1 \in \mathcal{P}(A_1)$ . Then  $|P'_1| + \varepsilon > l(A_1)$ ,  $|P'_2| + \varepsilon > l(A_2)$ , and both  $P'_1$  and  $P'_2$  determine the same sequence of angles. Using Corollary 3.4.5

$$\frac{|P'_1|}{|P'_2|} = \frac{R}{r}$$

and so

$$l(A_2) < |P'_2| + \varepsilon = \frac{r}{R} |P'_1| + \varepsilon \leq \frac{r}{R} l(A_1) + \varepsilon.$$

Since  $\varepsilon$  is arbitrary, this shows  $Rl(A_2) \leq rl(A_1)$ . But now reverse the argument and write

$$l(A_1) < |P'_1| + \varepsilon = \frac{R}{r} |P'_2| + \varepsilon \leq \frac{R}{r} l(A_2) + \varepsilon$$

which implies, since  $\varepsilon$  is arbitrary that  $Rl(A_2) \geq rl(A_1)$  and this has proved the following theorem.

**Theorem 3.5.1** *Let  $\theta$  be an angle which subtends two arcs,  $A_R$  on a circle of radius  $R$  and  $A_r$  on a circle of radius  $r$ . Then denoting by  $l(A)$  the length of a circular arc as described above,  $Rl(A_r) = rl(A_R)$ .*

Before proceeding further, note the proof of the above theorem involved showing  $l(A_1) < \frac{R}{r}l(A_2) + \varepsilon$  where  $\varepsilon > 0$  was arbitrary and from this, the conclusion that  $l(A_1) \leq \frac{R}{r}l(A_2)$ . This is a very typical way of showing one number is no larger than another. To show  $a \leq b$  first show that for every  $\varepsilon > 0$  it follows that  $a < b + \varepsilon$ . This implies  $a - b < \varepsilon$  for all positive  $\varepsilon$  and so it must be the case that  $a - b \leq 0$  since otherwise, you could take  $\varepsilon = \frac{a-b}{2}$  and conclude  $0 < a - b < \frac{a-b}{2}$ , a contradiction.

With this preparation, here is the definition of the measure of an angle.

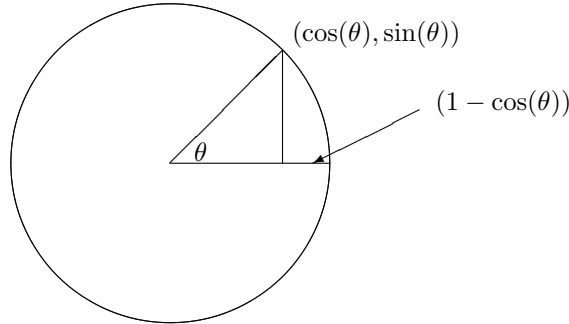
**Definition 3.5.2** *Let  $\theta$  be an angle. The measure of  $\theta$  is defined to be the length of the circular arc subtended by  $\theta$  on a circle of radius  $r$  divided by  $r$ . This is also called the radian measure of the angle.*

You should note the measure of  $\theta$  is independent of dimension. This is because the units of length cancel when the division takes place.

**Proposition 3.5.3** *The above definition is well defined and also, if  $A$  is an arc subtended by the angle  $\theta$  on a circle of radius  $r$  then the length of  $A$ , denoted by  $l(A)$  is given by  $l(A) = r\theta$ .*

**Proof:** That the definition is well defined follows from Theorem 3.5.1. The formula also follows from Theorem 3.5.1 and letting  $R = 1$ .

Now is a good time to present a useful inequality which may or may not be self evident. Here is a picture which illustrates the conclusion of this corollary.



The following corollary states that the length of the subtended arc shown in the picture is longer than the vertical side of the triangle and smaller than the sum of the vertical side with the segment having length  $1 - \cos \theta$ . To me, this seems abundantly clear but in case it is hard to believe, the following corollary gives a proof.

**Corollary 3.5.4** *Let  $0 < \text{radian measure of } \theta < \pi/4$ . Then letting  $A$  be the arc on the unit circle resulting from situating the angle with one side on the positive  $x$  axis and the other side pointing up from the positive  $x$  axis,*

$$(1 - \cos \theta) + \sin \theta \geq l(A) \geq \sin \theta \quad (3.2)$$

**Proof:** Situate the angle,  $\theta$  such that one side is on the positive  $x$  axis and extend the other side till it intersects the unit circle at the point,  $(\cos \theta, \sin \theta)$ . Then denoting the resulting arc on the circle by  $A$ , it follows that for all  $P \in \mathcal{P}(A)$  the inequality  $(1 - \cos \theta) + \sin \theta \geq |P| \geq \sin \theta$ . It follows that  $(1 - \cos \theta) + \sin \theta$  is an upper bound for all the  $|P|$  where  $P \in \mathcal{P}(A)$  and so  $(1 - \cos \theta) + \sin \theta$  is at least as large as the sup or least upper bound of the

$|P|$ . This proves the top half of the inequality. The bottom half follows because  $l(A) \geq L$  where  $L$  is the length of the line segment joining  $(\cos \theta, \sin \theta)$  and  $(1, 0)$  due to the definition of  $l(A)$ . However,  $L \geq \sin \theta$  because  $L$  is the length of the hypotenuse of a right triangle having  $\sin \theta$  as one of the sides.

### 3.6 The Trigonometric Functions

Now the Trigonometric functions will be defined as functions of an arbitrary real variable. Up till now these have been defined as functions of pointy things called angles. The following theorem will make possible the definition.

**Theorem 3.6.1** *Let  $b \in \mathbb{R}$ . Then there exists a unique integer  $p$  and real number  $r$  such that  $0 \leq r < 2\pi$  and  $b = p2\pi + r$ .*

**Proof:** First suppose  $b \geq 0$ . Then from Theorem 2.8.11 on Page 29 there exists a unique integer,  $p$  such that  $b = 2\pi p + r$  where  $0 \leq r < 2\pi$ . Now suppose  $b < 0$ . Then there exists a unique integer,  $p$  such that  $-b = 2\pi p + r_1$  where  $r_1 \in [0, 2\pi)$ . If  $r_1 = 0$ , then  $b = (-p)2\pi$ .

Otherwise,  $b = (-p)2\pi + (-r_1) = (-p-1)2\pi + \left(\overbrace{2\pi - r_1}^{\equiv r}\right)$  and  $r \equiv 2\pi - r_1 \in (0, 2\pi)$ .

The following definition is for  $\sin b$  and  $\cos b$  for  $b \in \mathbb{R}$ .

**Definition 3.6.2** *Let  $b \in \mathbb{R}$ . Then  $\sin b \equiv \sin r$  and  $\cos b \equiv \cos r$  where  $b = 2\pi p + r$  for  $p$  an integer, and  $r \in [0, 2\pi)$ .*

Several observations are now obvious from this.

**Observation 3.6.3** *Let  $k \in \mathbb{Z}$ , then the following formulas hold.*

$$\sin b = -\sin(-b), \cos b = \cos(-b), \quad (3.3)$$

$$\sin(b + 2k\pi) = \sin b, \cos(b + 2k\pi) = \cos b \quad (3.4)$$

$$\cos^2 b + \sin^2 b = 1 \quad (3.5)$$

The other trigonometric functions are defined in the usual way as in (3.1) provided they make sense.

From the observation that the  $x$  and  $y$  axes intersect at right angles the four arcs on the unit circle subtended by these axes are all of equal length. Therefore, the measure of a right angle must be  $2\pi/4 = \pi/2$ . The measure of the angle which is determined by the arc from  $(1, 0)$  to  $(-1, 0)$  is seen to equal  $\pi$  by the same reasoning. From the definition of the trig functions,  $\cos(\pi/2) = 0$  and  $\sin(\pi/2) = 1$ . You can easily find other values for  $\cos$  and  $\sin$  at all the other multiples of  $\pi/2$ .

The next topic is the important formulas for the trig. functions of sums and differences of numbers. For  $b \in \mathbb{R}$ , denote by  $r_b$  the element of  $[0, 2\pi)$  having the property that  $b = 2\pi p + r_b$  for  $p$  an integer.

**Lemma 3.6.4** *Let  $x, y \in \mathbb{R}$ . Then  $r_{x+y} = r_x + r_y + 2k\pi$  for some  $k \in \mathbb{Z}$ .*

**Proof:** By definition,

$$x + y = 2\pi p + r_{x+y}, \quad x = 2\pi p_1 + r_x, \quad y = 2\pi p_2 + r_y.$$

From this the result follows because

$$0 = ((x+y) - x) - y = 2\pi \overbrace{((p - p_1) - p_2)}^{\equiv -k} + r_{x+y} - (r_x + r_y).$$

Let  $z \in \mathbb{R}$  and let  $p(z)$  denote the point on the unit circle determined by the length  $r_z$  whose coordinates are  $\cos z$  and  $\sin z$ . Thus, starting at  $(1, 0)$  and moving counter clockwise a distance of  $r_z$  on the unit circle yields  $p(z)$ . Note also that  $p(z) = p(r_z)$ .

**Lemma 3.6.5** *Let  $x, y \in \mathbb{R}$ . Then the length of the arc between  $p(x+y)$  and  $p(x)$  is equal to the length of the arc between  $p(y)$  and  $(1, 0)$ .*

**Proof:** The length of the arc between  $p(x+y)$  and  $p(x)$  is  $|r_{x+y} - r_x|$ . There are two cases to consider here.

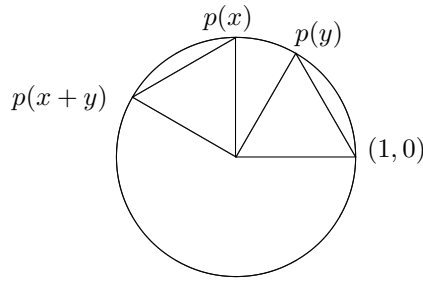
First assume  $r_{x+y} \geq r_x$ . Then  $|r_{x+y} - r_x| = r_{x+y} - r_x = r_y + 2k\pi$ . Since both  $r_{x+y}$  and  $r_x$  are in  $[0, 2\pi)$ , their difference is also in  $[0, 2\pi)$  and so  $k = 0$ . Therefore, the arc joining  $p(x)$  and  $p(x+y)$  is of the same length as the arc joining  $p(y)$  and  $(1, 0)$ . In the other case,  $r_{x+y} < r_x$  and in this case  $|r_{x+y} - r_x| = r_x - r_{x+y} = -r_y - 2k\pi$ . Since  $r_x$  and  $r_{x+y}$  are both in  $[0, 2\pi)$  their difference is also in  $[0, 2\pi)$  and so in this case  $k = -1$ . Therefore, in this case,  $|r_{x+y} - r_x| = 2\pi - r_y$ . Now since the circumference of the unit circle is  $2\pi$ , the length of the arc joining  $p(2\pi - r_y)$  to  $(1, 0)$  is the same as the length of the arc joining  $p(r_y) = p(y)$  to  $(1, 0)$ . This proves the lemma.

The following theorem is the fundamental identity from which all the major trig. identities involving sums and differences of angles are derived.

**Theorem 3.6.6** *Let  $x, y \in \mathbb{R}$ . Then*

$$\cos(x+y)\cos y + \sin(x+y)\sin y = \cos x. \quad (3.6)$$

**Proof:** Recall that for a real number,  $z$ , there is a unique point,  $p(z)$  on the unit circle and the coordinates of this point are  $\cos z$  and  $\sin z$ . Now from the above lemma, the length of the arc between  $p(x+y)$  and  $p(x)$  has the same length as the arc between  $p(y)$  and  $p(0)$ . As an illustration see the following picture.



It follows from the definition of the radian measure of an angle that the two angles determined by these arcs are equal and so, by Corollary 3.4.5 the distance between the points  $p(x+y)$  and  $p(x)$  must be the same as the distance from  $p(y)$  to  $p(0)$ . Writing this in terms of the definition of the trig functions and the distance formula,

$$(\cos(x+y) - \cos x)^2 + (\sin(x+y) - \sin x)^2 = (\cos y - 1)^2 + \sin^2 y.$$

$$\begin{aligned}\cos^2(x+y) + \cos^2 x - 2\cos(x+y)\cos x + \sin^2(x+y) + \sin^2 x - 2\sin(x+y)\sin x \\ = \cos^2 y - 2\cos y + 1 + \sin^2 y\end{aligned}$$

From Observation 3.6.3 this implies (3.6). This proves the theorem.

Letting  $y = \pi/2$ , this shows that

$$\sin(x + \pi/2) = \cos x. \quad (3.7)$$

Now let  $u = x + y$  and  $v = y$ . Then (3.6) implies

$$\cos u \cos v + \sin u \sin v = \cos(u - v) \quad (3.8)$$

Also, from this and (3.3),

$$\begin{aligned}\cos(u + v) &= \cos(u - (-v)) \\ &= \cos u \cos(-v) + \sin u \sin(-v) \\ &= \cos u \cos v - \sin u \sin v\end{aligned} \quad (3.9)$$

Thus, letting  $v = \pi/2$ ,

$$\cos\left(u + \frac{\pi}{2}\right) = -\sin u. \quad (3.10)$$

It follows

$$\begin{aligned}\sin(x + y) &= -\cos\left(x + \frac{\pi}{2} + y\right) \\ &= -\left[\cos\left(x + \frac{\pi}{2}\right)\cos y - \sin\left(x + \frac{\pi}{2}\right)\sin y\right] \\ &= \sin x \cos y + \sin y \cos x\end{aligned} \quad (3.11)$$

Then using Observation 3.6.3 again, this implies

$$\sin(x - y) = \sin x \cos y - \cos x \sin y. \quad (3.12)$$

In addition to this, Observation 3.6.3 implies

$$\cos 2x = \cos^2 x - \sin^2 x \quad (3.13)$$

$$= 2\cos^2 x - 1 \quad (3.14)$$

$$= 1 - 2\sin^2 x \quad (3.15)$$

Therefore, making use of the above identities, and Observation 3.6.3,

$$\begin{aligned}\cos(3x) &= \cos 2x \cos x - \sin 2x \sin x \\ &= (2\cos^2 x - 1)\cos x - 2\cos x \sin^2 x \\ &= 4\cos^3 x - 3\cos x\end{aligned} \quad (3.16)$$

With these fundamental identities, it is easy to obtain the cosine and sine of many special angles, called reference angles. First,  $\cos(\frac{\pi}{4})$ .

$$0 = \cos\left(\frac{\pi}{2}\right) = \cos\left(\frac{\pi}{4} + \frac{\pi}{4}\right) = 2\cos^2\left(\frac{\pi}{4}\right) - 1$$

and so  $\cos(\frac{\pi}{4}) = \sqrt{2}/2$ . (Why do isn't it equal to  $-\sqrt{2}/2$ ? **Hint:** Draw a picture.) Thus  $\sin(\frac{\pi}{4}) = \sqrt{2}/2$  also. (Why?) Here is another one. From (3.16),

$$\begin{aligned}0 &= \cos\left(\frac{\pi}{2}\right) = \cos 3\left(\frac{\pi}{6}\right) \\ &= 4\cos^3\left(\frac{\pi}{6}\right) - 3\cos\left(\frac{\pi}{6}\right).\end{aligned}$$

Therefore,  $\cos\left(\frac{\pi}{6}\right) = \frac{\sqrt{3}}{2}$  and consequently,  $\sin\left(\frac{\pi}{6}\right) = \frac{1}{2}$ . Here is a short table including these and a few others. You should make sure you can obtain all these entries. In the table,  $\theta$  refers to the radian measure of the angle. From now on, angles are considered as real numbers, not as pointy things.

$\theta$	0	$\frac{\pi}{6}$	$\frac{\pi}{4}$	$\frac{\pi}{3}$	$\frac{\pi}{2}$
$\cos \theta$	1	$\frac{\sqrt{3}}{2}$	$\frac{\sqrt{2}}{2}$	$\frac{1}{2}$	0
$\sin \theta$	0	$\frac{1}{2}$	$\frac{\sqrt{2}}{2}$	$\frac{\sqrt{3}}{2}$	1

### 3.7 Exercises

- Find  $\cos \theta$  and  $\sin \theta$  for  $\theta \in \left\{\frac{2\pi}{3}, \frac{3\pi}{4}, \frac{5\pi}{6}, \pi, \frac{7\pi}{6}, \frac{5\pi}{4}, \frac{4\pi}{3}, \frac{3\pi}{2}, \frac{5\pi}{3}, \frac{7\pi}{4}, \frac{11\pi}{6}, 2\pi\right\}$ .
- Prove  $\cos^2 \theta = \frac{1+\cos 2\theta}{2}$  and  $\sin^2 \theta = \frac{1-\cos 2\theta}{2}$ .
- $\pi/12 = \pi/3 - \pi/4$ . Therefore, from Problem 2,  $\cos(\pi/12) = \sqrt{\frac{1+(\sqrt{3}/2)}{2}}$ . On the other hand,  

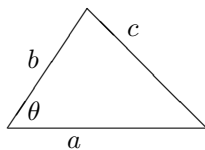
$$\cos(\pi/12) = \cos(\pi/3 - \pi/4) = \cos \pi/3 \cos \pi/4 + \sin \pi/3 \sin \pi/4$$
and so  $\cos(\pi/12) = \sqrt{2}/4 + \sqrt{6}/4$ . Is there a problem here? Please explain.
- Prove  $1 + \tan^2 \theta = \sec^2 \theta$  and  $1 + \cot^2 \theta = \csc^2 \theta$ .
- Prove that  $\sin x \cos y = \frac{1}{2}(\sin(x+y) + \sin(x-y))$ .
- Prove that  $\sin x \sin y = \frac{1}{2}(\cos(x-y) - \cos(x+y))$ .
- Prove that  $\cos x \cos y = \frac{1}{2}(\cos(x+y) + \cos(x-y))$ .
- Using Problem 5, find an identity for  $\sin x - \sin y$ .
- Suppose  $\sin x = a$  where  $0 < a < 1$ . Find all possible values for
  - $\tan x$
  - $\cot x$
  - $\sec x$
  - $\csc x$
  - $\cos x$
- Solve the equations and give all solutions.
  - $\sin(3x) = \frac{1}{2}$
  - $\cos(5x) = \frac{\sqrt{3}}{2}$
  - $\tan(x) = \sqrt{3}$
  - $\sec(x) = 2$
  - $\sin(x+7) = \frac{\sqrt{2}}{2}$
  - $\cos^2(x) = \frac{1}{2}$
  - $\sin^4(x) = 4$
- Find a formula for  $\tan(x+y)$  in terms of trig. functions of  $x$  and  $y$ .

12. Find a formula for  $\tan(2x)$  in terms of trig. functions of  $x$ .
13. Find a formula for  $\tan\left(\frac{x}{2}\right)$  in terms of trig. functions of  $x$ .
14. Sketch a graph of  $y = \sin x$ .
15. Sketch a graph of  $y = \cos x$ .
16. Sketch a graph of  $y = \sin 2x$ .
17. Sketch a graph of  $y = \tan x$ .
18. Using Problem 2 graph  $y = \cos^2 x$ .
19. If  $f(x) = A \cos \alpha x + B \sin \alpha x$ , show there exists  $\phi$  such that

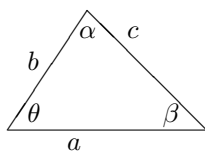
$$f(x) = \sqrt{A^2 + B^2} \sin(\alpha x + \phi).$$

Show there also exists  $\psi$  such that  $f(x) = \sqrt{A^2 + B^2} \cos(\alpha x + \psi)$ . This is a very important result, enough that some of these quantities are given names.  $\sqrt{A^2 + B^2}$  is called the amplitude and  $\phi$  or  $\psi$  are called phase shifts.

20. Using Problem 19 graph  $y = \sin x + \sqrt{3} \cos x$ .
21. Give all solutions to  $\sin x + \sqrt{3} \cos x = \sqrt{3}$ . **Hint:** Use Problem 20.
22. If  $ABC$  is a triangle where the capitol letters denote vertices of the triangle and the angle at the vertex. Let  $a$  be the length of the side opposite  $A$  and  $b$  is the length of the side opposite  $B$  and  $c$  is the length of the side opposite the vertex,  $C$ . The law of sines says  $\sin(A)/a = \sin(B)/b = \sin(C)/c$ . Prove the law of sines from the definition of the trigonometric functions.
23. In the picture,  $a = 5, b = 3$ , and  $\theta = \frac{2}{3}\pi$ . Find  $c$ .



24. In the picture,  $\theta = \frac{1}{4}\pi$ ,  $\alpha = \frac{2}{3}\pi$  and  $c = 3$ . Find  $a$ .

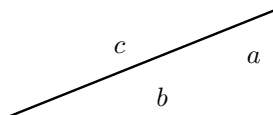




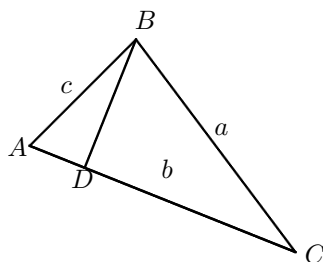
## 3.8 Some Basic Area Formulas

### 3.8.1 Areas Of Triangles And Parallelograms

This section is a review of how to find areas of some simple figures. The discussion will be somewhat informal since it is assumed the reader has seen this sort of thing already. First of all, consider a right triangle as indicated in the following picture.

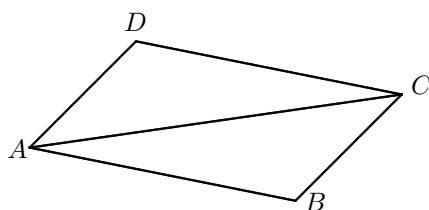


The area of this triangle shown above must equal  $ab/2$  because it is half of a rectangle having sides  $a$  and  $b$ . Now consider a general triangle in which a line perpendicular to the line from  $A$  to  $C$  has been drawn through  $B$ .



The area of this triangle would be the sum of the two right triangles formed. Thus this area would be  $\frac{1}{2} (\overline{BD}) (\overline{AD} + \overline{CD}) = \frac{1}{2} (\overline{BD}) b$ . In words, the area of the triangle equals one half the base times the height. This also holds if the height and base are chosen with respect to any other side of the triangle.

A parallelogram is a four sided figure which is formed when two identical triangles are joined along a corresponding side with the corresponding angles not adjacent. For example, see the picture in which the two triangles are  $ABC$  and  $CDA$ .



Note the height of triangle  $ABC$  taken with respect to side  $AB$  is the same as the height of the parallelogram taken with respect to this same side. Therefore, the area of this parallelogram equals twice the area of one of these triangles which equals  $2\overline{AB}$  (height of parallelogram)  $\frac{1}{2} = \overline{AB}$  (height of parallelogram). Similarly the area equals height times base where the base is any side of the parallelogram and the height is taken with respect to that side, as just described in the case where  $AB$  is the side.

### 3.8.2 The Area Of A Circular Sector

Consider an arc,  $A$ , of a circle of radius  $r$  which subtends an angle,  $\theta$ . The circular sector determined by  $A$  is obtained by joining the ends of the arc,  $A$ , to the center of the circle. The sector,  $S(\theta)$  denotes the points which lie between the arc,  $A$  and the two lines just mentioned. The angle between the two lines is called the central angle of the sector. The problem is to define the area of this shape. First a fundamental inequality must be obtained.

**Lemma 3.8.1** *Let  $1 > \varepsilon > 0$  be given. Then whenever the positive number,  $\alpha$ , is small enough,*

$$1 \leq \frac{\alpha}{\sin \alpha} \leq 1 + \varepsilon \quad (3.17)$$

and

$$1 + \varepsilon \geq \frac{\alpha}{\tan \alpha} \geq 1 - \varepsilon \quad (3.18)$$

**Proof:** This follows from Corollary 3.5.4 on Page 59. In this corollary,  $l(A) = \alpha$  and so

$$1 - \cos \alpha + \sin \alpha \geq \alpha \geq \sin \alpha.$$

Therefore, dividing by  $\sin \alpha$ ,

$$\frac{1 - \cos \alpha}{\sin \alpha} + 1 \geq \frac{\alpha}{\sin \alpha} \geq 1. \quad (3.19)$$

Now using the properties of the trig functions,

$$\begin{aligned} \frac{1 - \cos \alpha}{\sin \alpha} &= \frac{1 - \cos^2 \alpha}{\sin \alpha (1 + \cos \alpha)} \\ &= \frac{\sin^2 \alpha}{\sin \alpha (1 + \cos \alpha)} = \frac{\sin \alpha}{1 + \cos \alpha}. \end{aligned}$$

From the definition of the sin and cos, whenever  $\alpha$  is small enough,

$$\frac{\sin \alpha}{1 + \cos \alpha} < \varepsilon$$

and so (3.19) implies that for such  $\alpha$ , (3.17) holds. To obtain (3.18), let  $\alpha$  be small enough that (3.17) holds and multiply by  $\cos \alpha$ . Then for such  $\alpha$ ,

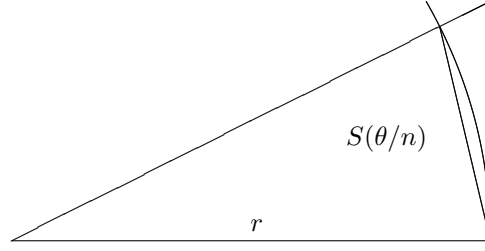
$$\cos \alpha \leq \frac{\alpha}{\tan \alpha} \leq (1 + \varepsilon) \cos \alpha$$

Taking  $\alpha$  smaller if necessary and noting that for all  $\alpha$  small enough,  $\cos \alpha$  is very close to 1, yields (3.18). This proves the lemma.

This lemma is very important in another context.

**Theorem 3.8.2** *Let  $S(\theta)$  denote the sector of a circle of radius  $r$  having central angle,  $\theta$ . Then the area of  $S(\theta)$  equals  $\frac{r^2}{2}\theta$ .*

**Proof:** Let the angle which  $A$  subtends be denoted by  $\theta$  and divide this sector into  $n$  equal sectors each of which has a central angle equal to  $\theta/n$ . The following is a picture of one of these.



In the picture, there is a circular sector,  $S(\theta/n)$  and inside this circular sector is a triangle while outside the circular sector is another triangle. Thus any reasonable definition of area would require

$$\frac{r^2}{2} \sin(\theta/n) \leq \text{area of } S(\theta/n) \leq \frac{r^2}{2} \tan(\theta/n).$$

It follows the area of the whole sector having central angle  $\theta$  must satisfy the following inequality.

$$\frac{nr^2}{2} \sin(\theta/n) \leq \text{area of } S(\theta) \leq \frac{nr^2}{2} \tan(\theta/n).$$

Therefore, for all  $n$ , the area of  $S(\theta)$  is trapped between the two numbers,

$$\frac{r^2}{2} \theta \frac{\sin(\theta/n)}{(\theta/n)}, \quad \frac{r^2}{2} \theta \frac{\tan(\theta/n)}{(\theta/n)}.$$

Now let  $\varepsilon > 0$  be given, a small positive number less than 1, and let  $n$  be large enough that

$$1 \geq \frac{\sin(\theta/n)}{(\theta/n)} \geq \frac{1}{1+\varepsilon}$$

and

$$\frac{1}{1+\varepsilon} \leq \frac{\tan(\theta/n)}{(\theta/n)} \leq \frac{1}{1-\varepsilon}.$$

Therefore,

$$\frac{r^2}{2} \theta \left( \frac{1}{1+\varepsilon} \right) \leq \text{Area of } S(\theta) \leq \left( \frac{1}{1-\varepsilon} \right) \frac{r^2}{2} \theta.$$

Since  $\varepsilon$  is an arbitrary small positive number, it follows the area of the sector equals  $\frac{r^2}{2} \theta$  as claimed. (Why?)

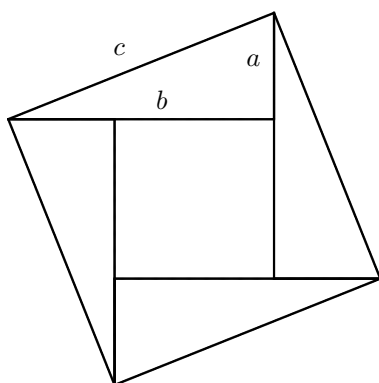
### 3.9 Exercises

1. Give another argument which verifies the Pythagorean theorem by supplying the details for the following argument<sup>3</sup>. Take the given right triangle and situate copies of it as shown below. The big four sided figure which results is a rectangle because all the angles are equal. Now from the picture, the area of the big square equals  $c^2$ , the area of each triangle equals  $ab/2$ , since it is half of a rectangle of area  $ab$ , and the area

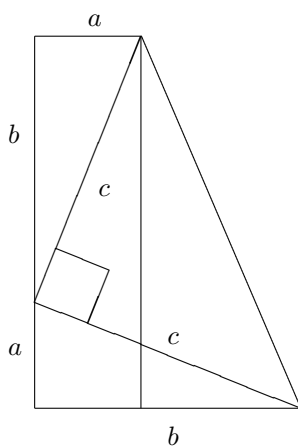
<sup>3</sup>This argument is old and was known to the Indian mathematician Bhaskar who lived 1114-1185 A.D.

of the inside square equals  $(b - a)^2$ . Here  $a$ ,  $b$ , and  $c$  are the lengths of the respective sides. Therefore,

$$\begin{aligned} c^2 &= 4(ab/2) + (b - a)^2 \\ &= 2ab + b^2 + a^2 - 2ab \\ &= a^2 + b^2. \end{aligned}$$



2. Another very simple and convincing proof of the Pythagorean theorem<sup>4</sup> is based on writing the area of the following trapezoid two ways. Sum the areas of three triangles in the following picture or write the area of the trapezoid as  $(a + b)a + \frac{1}{2}(a + b)(b - a)$  which is the sum of a triangle and a rectangle as shown. Do it both ways and see the pythagorean theorem appear.



3. A right circular cone has radius  $r$  and height  $h$ . This is like an ice cream cone. Find the area of the side of this cone in terms of  $h$  and  $r$ . **Hint:** Think of painting the side of the cone and while the paint is still wet, rolling it on the floor yielding a circular sector.

---

<sup>4</sup>This argument involving the area of a trapezoid is due to James Garfield who was one of the presidents of the United States.

4. An equilateral triangle is one in which all sides are of equal length. Find the area of an equilateral triangle whose sides have length  $l$ .
5. Draw two parallel lines one having length  $a$  and the other having length  $b$  suppose also these lines are at a distance of  $h$  from each other. Now join the ends of these lines to obtain a four sided figure. What is the area of this four sided figure?

## 3.10 Parabolas, Ellipses, and Hyperbolas

### 3.10.1 The Parabola

A parabola is a collection of points,  $P$  in the plane such that the distance from  $P$  to a fixed line is the same as the distance from  $P$  to a given point,  $P_0$ . From this definition, one can obtain an equation which will describe a parabola. Suppose then that the line is  $y = c$  and the point is  $(a, b)$  where  $b \neq c$  as shown in the picture.

$$\underline{y = c}$$

$$P_0 = (a, b)$$

The distance from the point,  $P = (x, y)$  to the line is  $|c - y|$ . Therefore, the description of the parabola requires that

$$\sqrt{(x - a)^2 + (y - b)^2} = |c - y|.$$

Squaring both sides,

$$x^2 - 2xa + a^2 + y^2 - 2yb + b^2 = c^2 - 2cy + y^2$$

and so

$$(x - a)^2 + b^2 - c^2 = (2b - 2c)y. \quad (3.20)$$

The simplest case is when  $a = 0$  and  $b = -c$ . Then in this case, it reduces to

$$x^2 = 4cy.$$

Now consider an arbitrary equation of the form,  $y = dx^2 + ex + f$  where  $d \neq 0$ . By this is meant the set of points  $(x, y)$  such that the equation holds. Such a set of points always is a parabola. To see this, complete the square on the right as follows:

$$\begin{aligned} y &= dx^2 + ex + f \\ &= d \left( x^2 + \frac{e}{d}x + \frac{f}{d} \right) \\ &= d \left( x^2 + \frac{e}{d}x + \frac{e^2}{4d^2} \right) - \frac{e^2d}{4d^2} \\ &= d \left( x - \left( \frac{-e}{2d} \right) \right)^2 - \frac{e^2d}{4d^2}. \end{aligned}$$

Therefore, letting  $a = \frac{-e}{2d}$ ,

$$\frac{1}{d}y = (x - a)^2 - \frac{e^2}{4d^2}$$

Now you can show that there exists numbers,  $b$  and  $c$  such that

$$-\frac{e^2}{4d^2} = b^2 - c^2, \quad \frac{1}{d} = 2b - 2c. \quad (3.21)$$

Then the above equation reduces to (3.20).

The line,  $y = c$  is called the directrix and the point,  $P_0$  in the above is called the focus. Exactly similar results occur if the directrix is of the form  $x = c$  and by similar arguments to those above, the set of points in the plane satisfying  $ay^2 + by + c = x$  is also a parabola.

### 3.10.2 The Ellipse

With an ellipse, there are two points,  $P_1$  and  $P_2$  which are fixed and the ellipse consists of the set of points,  $P$  such that  $d(P, P_1) + d(P, P_2) = c$ , where  $c$  is a fixed positive number. These two points are called the foci of the ellipse. Each is called a focus point by itself. Now one can obtain an equation which will describe an ellipse much as was done with the parabola. Let the two given points be  $(a, b)$  and  $(a, b + h)$ . Let a generic point on the ellipse be  $(x, y)$ . Then according to the description of an ellipse and the distance formula,

$$\sqrt{(x-a)^2 + (y-b)^2} + \sqrt{(x-a)^2 + (y-b-h)^2} = c. \quad (3.22)$$

Subtracting  $\sqrt{(x-a)^2 + (y-b)^2}$  from both sides,

$$\sqrt{(x-a)^2 + (y-b-h)^2} = c - \sqrt{(x-a)^2 + (y-b)^2} \geq 0.$$

Now squaring both sides yields

$$\begin{aligned} (x-a)^2 + (y-b-h)^2 &= c^2 - 2\sqrt{(x-a)^2 + (y-b)^2}c \\ &\quad + (x-a)^2 + (y-b)^2. \end{aligned}$$

Therefore,

$$\begin{aligned} (y-b)^2 - 2h(y-b) + h^2 &= c^2 - 2\sqrt{(x-a)^2 + (y-b)^2}c \\ &\quad + (y-b)^2 \end{aligned}$$

and so

$$-2h(y-b) + h^2 = c^2 - 2\sqrt{(x-a)^2 + (y-b)^2}c$$

Therefore,

$$-2h(y-b) + h^2 - c^2 = -2\sqrt{(x-a)^2 + (y-b)^2}c$$

Now square both sides again to obtain

$$(h^2 - c^2)^2 - 4h(y-b)(h^2 - c^2) + 4h^2(y-b)^2 = 4c^2((x-a)^2 + (y-b)^2).$$

Simplifying this yields

$$(c^2 - h^2)(y-b)^2 + h(y-b)(h^2 - c^2) + c^2(x-a)^2 = \frac{1}{4}(c^2 - h^2)^2$$

which simplifies further to

$$(y-b)^2 - h(y-b) + \frac{h^2}{4} + \left(\frac{c^2}{c^2 - h^2}\right)(x-a)^2 = \frac{1+h^2}{4}$$

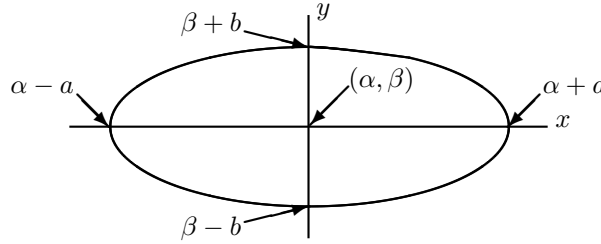
which is equivalent to

$$\frac{(y - b - \frac{h}{2})^2}{(\frac{1+h^2}{4})} + \frac{(x - a)^2}{(\frac{1+h^2}{4} / \frac{c^2}{c^2 - h^2})} = 1.$$

Thus, redefining the constants, an ellipse has the form

$$\frac{(y - \beta)^2}{b^2} + \frac{(x - \alpha)^2}{a^2} = 1. \quad (3.23)$$

Note that if  $a = b = r$ , this reduces to the equation for a circle of radius  $r$  centered at the point  $(\alpha, \beta)$ . This last expression is the generic equation for an ellipse. Here is the graph of a typical ellipse.



In this ellipse,  $b < a$ . If  $b > a$ , the ellipse would be long in the  $y$  direction rather than the  $x$  direction. (Why?) Suppose  $(x_1, y_1)$  and  $(x_2, y_2)$  are two points on the above ellipse in which  $b > a$ , then  $|x_1 - x_2| \leq 2a$  because from the above equation, it follows that  $|x_i - \alpha| \leq a$  for  $i = 1, 2$  implying that

$$|x_1 - x_2| \leq |x_1 - \alpha| + |\alpha - x_2| \leq a + a = 2a$$

and similarly,  $|y_1 - y_2| \leq 2b$ . Therefore,

$$|(x_1, y_1) - (x_2, y_2)| \leq \sqrt{4a^2 + 4b^2} \leq 2\sqrt{a^2 + b^2} \leq 2b$$

Thus the greatest distance between two points on the ellipse equals  $2b$  and occurs when the two points are  $(\alpha, b + \beta)$  and  $(\alpha, \beta - b)$ . This greatest distance between any two points is called the diameter and this shows the diameter of an ellipse is twice the larger of the two numbers appearing in the denominators on the left in (3.23).

### 3.10.3 The Hyperbola

With a hyperbola, there are two points,  $P_1$  and  $P_2$  which are fixed and the hyperbola consists of the set of points,  $P$  such that  $d(P, P_1) - d(P, P_2) = c$ , where  $c$  is a fixed positive number. These two points are called the foci of the hyperbola. Each is called a focus point by itself. Now one can obtain an equation which will describe a hyperbola. Let the two given points be  $(a, b)$  and  $(a, b + h)$ . Let a generic point on the hyperbola be  $(x, y)$ . Then according to the description of a hyperbola and the distance formula,

$$\sqrt{(x - a)^2 + (y - b)^2} - \sqrt{(x - a)^2 + (y - b - h)^2} = c.$$

You can now show that the equation of a hyperbola is of the form

$$\frac{(x - \alpha)^2}{a^2} - \frac{(y - \beta)^2}{b^2} = 1 \quad (3.24)$$

or

$$\frac{(y - \beta)^2}{b^2} - \frac{(x - \alpha)^2}{a^2} = 1. \quad (3.25)$$

If you like, you can simply take these last two equations as the definition of a hyperbola.

### 3.11 Exercises

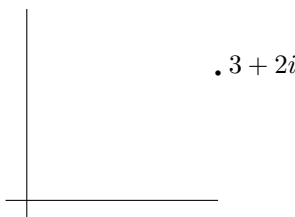
1. Consider  $y = 2x^2 + 3x + 7$ . Find the focus and the directrix of this parabola.
2. Find the numbers,  $b, c$  which make (3.21) hold.
3. Derive a similar formula to (3.20) in the case that the directrix is of the form  $x = c$ .
4. Sketch a graph of the ellipse whose equation is  $\frac{(x-1)^2}{4} + \frac{(y-2)^2}{9} = 1$ .
5. Sketch a graph of the ellipse whose equation is  $\frac{(x-1)^2}{9} + \frac{(y-2)^2}{4} = 1$ .
6. Sketch a graph of the hyperbola,  $\frac{x^2}{4} - \frac{y^2}{9} = 1$ .
7. Sketch a graph of the hyperbola,  $\frac{y^2}{4} - \frac{x^2}{9} = 1$ .
8. What is the diameter of the ellipse,  $\frac{(x-1)^2}{9} + \frac{(y-2)^2}{4} = 1$ .
9. Verify that either (3.24) or (3.25) holds for a hyperbola.
10. Show that the set of points which satisfies either (3.24) or (3.25) is unbounded. (If  $n$  is any positive number there exist points  $(x, y)$  satisfying the equations given such that  $|(x, y)| > n$ .)



# The Complex Numbers

This chapter gives a brief treatment of the complex numbers. This will not be needed in Calculus but you will need it when you take differential equations and various other subjects so it is a good idea to consider the subject. These things used to be taught in precalculus classes and people were expected to know them before taking calculus. However, if you are in a hurry to get to calculus, you can skip this short chapter.

Just as a real number should be considered as a point on the line, a complex number is considered a point in the plane which can be identified in the usual way using the Cartesian coordinates of the point. Thus  $(a, b)$  identifies a point whose  $x$  coordinate is  $a$  and whose  $y$  coordinate is  $b$ . In dealing with complex numbers, such a point is written as  $a + ib$ . For example, in the following picture, I have graphed the point  $3 + 2i$ . You see it corresponds to the point in the plane whose coordinates are  $(3, 2)$ .



Multiplication and addition are defined in the most obvious way subject to the convention that  $i^2 = -1$ . Thus,

$$(a + ib) + (c + id) = (a + c) + i(b + d)$$

and

$$\begin{aligned} (a + ib)(c + id) &= ac + iad + ibc + i^2bd \\ &= (ac - bd) + i(bc + ad). \end{aligned}$$

Every non zero complex number,  $a + ib$ , with  $a^2 + b^2 \neq 0$ , has a unique multiplicative inverse.

$$\frac{1}{a + ib} = \frac{a - ib}{a^2 + b^2} = \frac{a}{a^2 + b^2} - i \frac{b}{a^2 + b^2}.$$

You should prove the following theorem.

**Theorem 4.0.1** *The complex numbers with multiplication and addition defined as above form a field satisfying all the field axioms listed on Page 15.*

The field of complex numbers is denoted as  $\mathbb{C}$ . An important construction regarding complex numbers is the complex conjugate denoted by a horizontal line above the number.

It is defined as follows.

$$\overline{a + ib} \equiv a - ib.$$

What it does is reflect a given complex number across the  $x$  axis. Algebraically, the following formula is easy to obtain.

$$(\overline{a + ib})(a + ib) = a^2 + b^2.$$

**Definition 4.0.2** Define the absolute value of a complex number as follows.

$$|a + ib| \equiv \sqrt{a^2 + b^2}.$$

Thus, denoting by  $z$  the complex number,  $z = a + ib$ ,

$$|z| = (z\bar{z})^{1/2}.$$

With this definition, it is important to note the following. Be sure to verify this. It is not too hard but you need to do it.

**Remark 4.0.3** : Let  $z = a + ib$  and  $w = c + id$ . Then  $|z - w| = \sqrt{(a - c)^2 + (b - d)^2}$ . Thus the distance between the point in the plane determined by the ordered pair,  $(a, b)$  and the ordered pair  $(c, d)$  equals  $|z - w|$  where  $z$  and  $w$  are as just described.

For example, consider the distance between  $(2, 5)$  and  $(1, 8)$ . From the distance formula this distance equals  $\sqrt{(2 - 1)^2 + (5 - 8)^2} = \sqrt{10}$ . On the other hand, letting  $z = 2 + i5$  and  $w = 1 + i8$ ,  $z - w = 1 - i3$  and so  $(z - w)(\overline{z - w}) = (1 - i3)(1 + i3) = 10$  so  $|z - w| = \sqrt{10}$ , the same thing obtained with the distance formula.

Complex numbers, are often written in the so called polar form which is described next. Suppose  $x + iy$  is a complex number. Then

$$x + iy = \sqrt{x^2 + y^2} \left( \frac{x}{\sqrt{x^2 + y^2}} + i \frac{y}{\sqrt{x^2 + y^2}} \right).$$

Now note that

$$\left( \frac{x}{\sqrt{x^2 + y^2}} \right)^2 + \left( \frac{y}{\sqrt{x^2 + y^2}} \right)^2 = 1$$

and so

$$\left( \frac{x}{\sqrt{x^2 + y^2}}, \frac{y}{\sqrt{x^2 + y^2}} \right)$$

is a point on the unit circle. Therefore, there exists a unique angle,  $\theta \in [0, 2\pi)$  such that

$$\cos \theta = \frac{x}{\sqrt{x^2 + y^2}}, \sin \theta = \frac{y}{\sqrt{x^2 + y^2}}.$$

The polar form of the complex number is then

$$r(\cos \theta + i \sin \theta)$$

where  $\theta$  is this angle just described and  $r = \sqrt{x^2 + y^2}$ .

A fundamental identity is the formula of De Moivre which follows.

**Theorem 4.0.4** Let  $r > 0$  be given. Then if  $n$  is a positive integer,

$$[r(\cos t + i \sin t)]^n = r^n (\cos nt + i \sin nt).$$

**Proof:** It is clear the formula holds if  $n = 1$ . Suppose it is true for  $n$ .

$$[r(\cos t + i \sin t)]^{n+1} = [r(\cos t + i \sin t)]^n [r(\cos t + i \sin t)]$$

which by induction equals

$$\begin{aligned} &= r^{n+1} (\cos nt + i \sin nt) (\cos t + i \sin t) \\ &= r^{n+1} ((\cos nt \cos t - \sin nt \sin t) + i (\sin nt \cos t + \cos nt \sin t)) \\ &= r^{n+1} (\cos (n+1)t + i \sin (n+1)t) \end{aligned}$$

by the formulas for the cosine and sine of the sum of two angles.

**Corollary 4.0.5** *Let  $z$  be a non zero complex number. Then there are always exactly  $k$   $k^{\text{th}}$  roots of  $z$  in  $\mathbb{C}$ .*

**Proof:** Let  $z = x + iy$  and let  $z = |z|(\cos t + i \sin t)$  be the polar form of the complex number. By De Moivre's theorem, a complex number,

$$r(\cos \alpha + i \sin \alpha),$$

is a  $k^{\text{th}}$  root of  $z$  if and only if

$$r^k (\cos k\alpha + i \sin k\alpha) = |z| (\cos t + i \sin t).$$

This requires  $r^k = |z|$  and so  $r = |z|^{1/k}$  and also both  $\cos(k\alpha) = \cos t$  and  $\sin(k\alpha) = \sin t$ . This can only happen if

$$k\alpha = t + 2l\pi$$

for  $l$  an integer. Thus

$$\alpha = \frac{t + 2l\pi}{k}, l \in \mathbb{Z}$$

and so the  $k^{\text{th}}$  roots of  $z$  are of the form

$$|z|^{1/k} \left( \cos \left( \frac{t + 2l\pi}{k} \right) + i \sin \left( \frac{t + 2l\pi}{k} \right) \right), l \in \mathbb{Z}.$$

Since the cosine and sine are periodic of period  $2\pi$ , there are exactly  $k$  distinct numbers which result from this formula.

**Example 4.0.6** *Find the three cube roots of  $i$ .*

First note that  $i = 1 \left( \cos \left( \frac{\pi}{2} \right) + i \sin \left( \frac{\pi}{2} \right) \right)$ . Using the formula in the proof of the above corollary, the cube roots of  $i$  are

$$1 \left( \cos \left( \frac{(\pi/2) + 2l\pi}{3} \right) + i \sin \left( \frac{(\pi/2) + 2l\pi}{3} \right) \right)$$

where  $l = 0, 1, 2$ . Therefore, the roots are

$$\cos \left( \frac{\pi}{6} \right) + i \sin \left( \frac{\pi}{6} \right), \cos \left( \frac{5}{6}\pi \right) + i \sin \left( \frac{5}{6}\pi \right),$$

and

$$\cos \left( \frac{3}{2}\pi \right) + i \sin \left( \frac{3}{2}\pi \right).$$

Thus the cube roots of  $i$  are  $\frac{\sqrt{3}}{2} + i \left( \frac{1}{2} \right)$ ,  $\frac{-\sqrt{3}}{2} + i \left( \frac{1}{2} \right)$ , and  $-i$ .

The ability to find  $k^{\text{th}}$  roots can also be used to factor some polynomials.

**Example 4.0.7** Factor the polynomial  $x^3 - 27$ .

First find the cube roots of 27. By the above procedure using De Moivre's theorem, these cube roots are  $3, 3\left(\frac{-1}{2} + i\frac{\sqrt{3}}{2}\right)$ , and  $3\left(\frac{-1}{2} - i\frac{\sqrt{3}}{2}\right)$ . Therefore,  $x^3 - 27 =$

$$(x - 3) \left( x - 3 \left( \frac{-1}{2} + i\frac{\sqrt{3}}{2} \right) \right) \left( x - 3 \left( \frac{-1}{2} - i\frac{\sqrt{3}}{2} \right) \right).$$

Note also  $\left( x - 3 \left( \frac{-1}{2} + i\frac{\sqrt{3}}{2} \right) \right) \left( x - 3 \left( \frac{-1}{2} - i\frac{\sqrt{3}}{2} \right) \right) = x^2 + 3x + 9$  and so

$$x^3 - 27 = (x - 3)(x^2 + 3x + 9)$$

where the quadratic polynomial,  $x^2 + 3x + 9$  cannot be factored without using complex numbers.

## 4.1 Exercises

1. Let  $z = 5 + i9$ . Find  $z^{-1}$ .
2. Let  $z = 2 + i7$  and let  $w = 3 - i8$ . Find  $zw, z + w, z^2$ , and  $w/z$ .
3. Give the complete solution to  $x^4 + 16 = 0$ .
4. Graph the complex cube roots of 8 in the complex plane. Do the same for the four fourth roots of 16.
5. If  $z$  is a complex number, show there exists  $\omega$  a complex number with  $|\omega| = 1$  and  $\omega z = |z|$ .
6. De Moivre's theorem says  $[r(\cos t + i \sin t)]^n = r^n (\cos nt + i \sin nt)$  for  $n$  a positive integer. Does this formula continue to hold for all integers,  $n$ , even negative integers? Explain.
7. You already know formulas for  $\cos(x + y)$  and  $\sin(x + y)$  and these were used to prove De Moivre's theorem. Now using De Moivre's theorem, derive a formula for  $\sin(5x)$  and one for  $\cos(5x)$ . **Hint:** Use Problem 18 on Page 32 and if you like, you might use Pascal's triangle to construct the binomial coefficients.
8. If  $z$  and  $w$  are two complex numbers and the polar form of  $z$  involves the angle  $\theta$  while the polar form of  $w$  involves the angle  $\phi$ , show that in the polar form for  $zw$  the angle involved is  $\theta + \phi$ . Also, show that in the polar form of a complex number,  $z$ ,  $r = |z|$ .
9. Factor  $x^3 + 8$  as a product of linear factors.
10. Write  $x^3 + 27$  in the form  $(x + 3)(x^2 + ax + b)$  where  $x^2 + ax + b$  cannot be factored any more using only real numbers.
11. Completely factor  $x^4 + 16$  as a product of linear factors.
12. Factor  $x^4 + 16$  as the product of two quadratic polynomials each of which cannot be factored further without using complex numbers.

13. If  $z, w$  are complex numbers prove  $\overline{zw} = \overline{z}\overline{w}$  and then show by induction that  $\overline{z_1 \cdots z_m} = \overline{z_1} \cdots \overline{z_m}$ . Also verify that  $\overline{\sum_{k=1}^m z_k} = \sum_{k=1}^m \overline{z_k}$ . In words this says the conjugate of a product equals the product of the conjugates and the conjugate of a sum equals the sum of the conjugates.
14. Suppose  $p(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$  where all the  $a_k$  are real numbers. Suppose also that  $p(z) = 0$  for some  $z \in \mathbb{C}$ . Show it follows that  $p(\overline{z}) = 0$  also.
15. I claim that  $1 = -1$ . Here is why.

$$-1 = i^2 = \sqrt{-1}\sqrt{-1} = \sqrt{(-1)^2} = \sqrt{1} = 1.$$

This is clearly a remarkable result but is there something wrong with it? If so, what is wrong?

16. De Moivre's theorem is really a grand thing. I plan to use it now for rational exponents, not just integers.

$$1 = 1^{(1/4)} = (\cos 2\pi + i \sin 2\pi)^{1/4} = \cos(\pi/2) + i \sin(\pi/2) = i.$$

Therefore, squaring both sides it follows  $1 = -1$  as in the previous problem. What does this tell you about De Moivre's theorem? Is there a profound difference between raising numbers to integer powers and raising numbers to non integer powers?



**Part II**

**Functions Of One Variable**





# Functions

## 5.1 General Considerations

The concept of a function is that of something which gives a unique output for a given input.

**Definition 5.1.1** Consider two sets,  $D$  and  $R$  along with a rule which assigns a unique element of  $R$  to every element of  $D$ . This rule is called a function and it is denoted by a letter such as  $f$ . The symbol,  $D(f) = D$  is called the domain of  $f$ . The set  $R$ , also written  $R(f)$ , is called the range of  $f$ . The set of all elements of  $R$  which are of the form  $f(x)$  for some  $x \in D$  is often denoted by  $f(D)$ . When  $R = f(D)$ , the function,  $f$ , is said to be onto. It is common notation to write  $f : D(f) \rightarrow R$  to denote the situation just described in this definition where  $f$  is a function defined on  $D$  having values in  $R$ .

**Example 5.1.2** Consider the list of numbers,  $\{1, 2, 3, 4, 5, 6, 7\} \equiv D$ . Define a function which assigns an element of  $D$  to  $R \equiv \{2, 3, 4, 5, 6, 7, 8\}$  by  $f(x) \equiv x + 1$  for each  $x \in D$ .

In this example there was a clearly defined procedure which determined the function. However, sometimes there is no discernible procedure which yields a particular function.

**Example 5.1.3** Consider the ordered pairs,  $(1, 2), (2, -2), (8, 3), (7, 6)$  and let

$$D \equiv \{1, 2, 8, 7\},$$

the set of first entries in the given set of ordered pairs,  $R \equiv \{2, -2, 3, 6\}$ , the set of second entries, and let  $f(1) = 2, f(2) = -2, f(8) = 3$ , and  $f(7) = 6$ .

Sometimes functions are not given in terms of a formula. For example, consider the following function defined on the positive real numbers having the following definition.

**Example 5.1.4** For  $x \in \mathbb{R}$  define

$$f(x) = \begin{cases} \frac{1}{n} & \text{if } x = \frac{m}{n} \text{ in lowest terms for } m, n \in \mathbb{Z} \\ 0 & \text{if } x \text{ is not rational} \end{cases} \quad (5.1)$$

This is a very interesting function called the Dirichlet function. Note that it is not defined in a simple way from a formula.

**Example 5.1.5** Let  $D$  consist of the set of people who have lived on the earth except for Adam and for  $d \in D$ , let  $f(d) \equiv$  the biological father of  $d$ . Then  $f$  is a function.

**Example 5.1.6** Consider a weight which is suspended at one end of a spring which is attached at the other end to the ceiling. Suppose the weight has extended the spring so that the force exerted by the spring exactly balances the force resulting from the weight on the spring. Measure the displacement of the mass,  $x$ , from this point with the positive direction being up, and define a function as follows:  $x(t)$  will equal the displacement of the spring at time  $t$  given knowledge of the velocity of the weight and the displacement of the weight at some particular time.

**Example 5.1.7** Certain chemicals decay with time. Suppose  $A_0$  is the amount of chemical at some given time. Then you could let  $A(t)$  denote the amount of the chemical at time  $t$ .

These last two examples show how physical problems can result in functions. Examples 5.1.6 and 5.1.7 are considered later in the book and techniques for finding  $x(t)$  and  $A(t)$  from the given conditions are presented.

In this chapter the functions are defined on some subset of  $\mathbb{R}$  having values in  $\mathbb{R}$ . Later this will be generalized. When  $D(f)$  is not specified, it is understood to consist of every number for which  $f$  makes sense. The following definition gives several ways to make new functions from old ones.

**Definition 5.1.8** Let  $f, g$  be functions with values in  $\mathbb{R}$ . Let  $a, b$  be elements of  $\mathbb{R}$ . Then  $af + bg$  is the name of a function whose domain is  $D(f) \cap D(g)$  which is defined as

$$(af + bg)(x) = af(x) + bg(x).$$

The function,  $fg$  is the name of a function which is defined on  $D(f) \cap D(g)$  given by

$$(fg)(x) = f(x)g(x).$$

Similarly for  $k$  an integer,  $f^k$  is the name of a function defined as

$$f^k(x) = (f(x))^k$$

The function,  $f/g$  is the name of a function whose domain is

$$D(f) \cap \{x \in D(g) : g(x) \neq 0\}$$

defined as

$$(f/g)(x) = f(x)/g(x).$$

If  $f : D(f) \rightarrow X$  and  $g : D(g) \rightarrow Y$ , then  $g \circ f$  is the name of a function whose domain is

$$\{x \in D(f) : f(x) \in D(g)\}$$

which is defined as

$$g \circ f(x) \equiv g(f(x)).$$

This is called the composition of the two functions.

You should note that  $f(x)$  is not a function. It is the value of the function at the point,  $x$ . The name of the function is  $f$ . Nevertheless, people often write  $f(x)$  to denote a function and it doesn't cause too many problems in beginning courses. When this is done, the variable,  $x$  should be considered as a generic variable free to be anything in  $D(f)$ . I will use this slightly sloppy abuse of notation whenever convenient. Thus,  $x^2 + 4$  may mean the function,  $f$ , given by  $f(x) = x^2 + 4$ .

**Example 5.1.9** Let  $f(t) = t$  and  $g(t) = 1 + t$ . Then  $fg : \mathbb{R} \rightarrow \mathbb{R}$  is given by

$$fg(t) = t(1 + t) = t + t^2.$$

**Example 5.1.10** Let  $f(t) = 2t + 1$  and  $g(t) = \sqrt{1 + t}$ . Then

$$g \circ f(t) = \sqrt{1 + (2t + 1)} = \sqrt{2t + 2}$$

for  $t \geq -1$ . If  $t < -1$  the inside of the square root sign is negative so makes no sense. Therefore,  $g \circ f : \{t \in \mathbb{R} : t \geq -1\} \rightarrow \mathbb{R}$ .

Note that in this last example, it was necessary to fuss about the domain of  $g \circ f$  because  $g$  is only defined for certain values of  $t$ .

The concept of a one to one function is very important. This is discussed in the following definition.

**Definition 5.1.11** For any function,  $f : D(f) \subseteq X \rightarrow Y$ , define the following set known as the inverse image of  $y$ .

$$f^{-1}(y) \equiv \{x \in D(f) : f(x) = y\}.$$

There may be many elements in this set, but when there is always only one element in this set for all  $y \in f(D(f))$ , the function  $f$  is one to one sometimes written,  $1 - 1$ . Thus  $f$  is one to one,  $1 - 1$ , if whenever  $f(x) = f(x_1)$ , then  $x = x_1$ . If  $f$  is one to one, the inverse function,  $f^{-1}$  is defined on  $f(D(f))$  and  $f^{-1}(y) = x$  where  $f(x) = y$ . Thus from the definition,  $f^{-1}(f(x)) = x$  for all  $x \in D(f)$  and  $f(f^{-1}(y)) = y$  for all  $y \in f(D(f))$ . Defining  $\text{id}$  by  $\text{id}(z) \equiv z$  this says  $f \circ f^{-1} = \text{id}$  and  $f^{-1} \circ f = \text{id}$ .

Polynomials and rational functions are particularly easy functions to understand.

**Definition 5.1.12** A function  $f$  is a polynomial if

$$f(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$$

where the  $a_i$  are real numbers and  $n$  is a nonnegative integer. In this case the degree of the polynomial,  $f(x)$  is  $n$ . Thus the degree of a polynomial is the largest exponent appearing on the variable.

$f$  is a rational function if

$$f(x) = \frac{h(x)}{g(x)}$$

where  $h$  and  $g$  are polynomials.

For example,  $f(x) = 3x^5 + 9x^2 + 7x + 5$  is a polynomial of degree 5 and

$$\frac{3x^5 + 9x^2 + 7x + 5}{x^4 + 3x + x + 1}$$

is a rational function.

Note that in the case of a rational function, the domain of the function might not be all of  $\mathbb{R}$ . For example, if

$$f(x) = \frac{x^2 + 8}{x + 1},$$

the domain of  $f$  would be all real numbers not equal to  $-1$ .

Closely related to the definition of a function is the concept of the graph of a function.

**Definition 5.1.13** Given two sets,  $X$  and  $Y$ , the Cartesian product of the two sets, written as  $X \times Y$ , is assumed to be a set described as follows.

$$X \times Y = \{(x, y) : x \in X \text{ and } y \in Y\}.$$

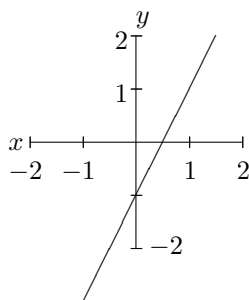
$\mathbb{R}^2$  denotes the Cartesian product of  $\mathbb{R}$  with  $\mathbb{R}$ .

The notion of Cartesian product is just an abstraction of the concept of identifying a point in the plane with an ordered pair of numbers.

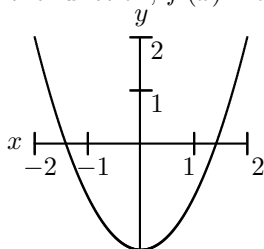
**Definition 5.1.14** Let  $f : D(f) \rightarrow R(f)$  be a function. The graph of  $f$  consists of the set,

$$\{(x, y) : y = f(x) \text{ for } x \in D(f)\}.$$

Note that knowledge of the graph of a function is equivalent to knowledge of the function. To find  $f(x)$ , simply observe the ordered pair which has  $x$  as its first element and the value of  $y$  equals  $f(x)$ . The graph of  $f$  can be represented by drawing a picture as mentioned earlier in the section on Cartesian coordinates beginning on Page 49. For example, consider the picture of a part of the graph of the function  $f(x) = 2x - 1$ .



Here is the graph of the function,  $f(x) = x^2 - 2$



**Definition 5.1.15** A function whose domain is defined as a set of the form  $\{k, k+1, k+2, \dots\}$  for  $k$  an integer is known as a sequence. Thus you can consider  $f(k)$ ,  $f(k+1)$ ,  $f(k+2)$ , etc. Usually the domain of the sequence is either  $\mathbb{N}$ , the natural numbers consisting of  $\{1, 2, 3, \dots\}$  or the nonnegative integers,  $\{0, 1, 2, 3, \dots\}$ . Also, it is traditional to write  $f_1, f_2$ , etc. instead of  $f(1)$ ,  $f(2)$ ,  $f(3)$  etc. when referring to sequences. In the above context,  $f_k$  is called the first term,  $f_{k+1}$  the second and so forth. It is also common to write the sequence, not as  $f$  but as  $\{f_i\}_{i=k}^{\infty}$  or just  $\{f_i\}$  for short.

**Example 5.1.16** Let  $\{a_k\}_{k=1}^{\infty}$  be defined by  $a_k \equiv k^2 + 1$ .

This gives a sequence. In fact,  $a_7 = a(7) = 7^2 + 1 = 50$  just from using the formula for the  $k^{\text{th}}$  term of the sequence.

It is nice when sequences come to us in this way from a formula for the  $k^{\text{th}}$  term. However, this is often not the case. Sometimes sequences are defined recursively. This happens, when the first several terms of the sequence are given and then a rule is specified which determines  $a_{n+1}$  from knowledge of  $a_1, \dots, a_n$ . This rule which specifies  $a_{n+1}$  from knowledge of  $a_k$  for  $k \leq n$  is known as a recurrence relation.

**Example 5.1.17** Let  $a_1 = 1$  and  $a_2 = 1$ . Assuming  $a_1, \dots, a_{n+1}$  are known,  $a_{n+2} \equiv a_n + a_{n+1}$ .

Thus the first several terms of this sequence, listed in order, are 1, 1, 2, 3, 5, 8,  $\dots$ . This particular sequence is called the Fibonacci sequence and is important in the study of reproducing rabbits.

**Definition 5.1.18** Let  $\{a_n\}$  be a sequence and let  $n_1 < n_2 < n_3, \dots$  be any strictly increasing list of integers such that  $n_1$  is at least as large as the first index used to define the sequence  $\{a_n\}$ . Then if  $b_k \equiv a_{n_k}$ ,  $\{b_k\}$  is called a subsequence of  $\{a_n\}$ .

For example, suppose  $a_n = (n^2 + 1)$ . Thus  $a_1 = 2$ ,  $a_3 = 10$ , etc. If

$$n_1 = 1, n_2 = 3, n_3 = 5, \dots, n_k = 2k - 1,$$

then letting  $b_k = a_{n_k}$ , it follows

$$b_k = ((2k - 1)^2 + 1) = 4k^2 - 4k + 2.$$

## 5.2 Exercises

- Let  $g(t) \equiv \sqrt{2-t}$  and let  $f(t) = \frac{1}{t}$ . Find  $g \circ f$ . Include the domain of  $g \circ f$ .
- Give the domains of the following functions.
  - $f(x) = \frac{x+3}{3x-2}$
  - $f(x) = \sqrt{x^2 - 4}$
  - $f(x) = \sqrt{4 - x^2}$
  - $f(x) = \sqrt{\frac{x-4}{3x+5}}$
  - $f(x) = \sqrt{\frac{x^2-4}{x+1}}$
- Let  $f: \mathbb{R} \rightarrow \mathbb{R}$  be defined by  $f(t) \equiv t^3 + 1$ . Is  $f$  one to one? Can you find a formula for  $f^{-1}$ ?
- Suppose  $a_1 = 1, a_2 = 3$ , and  $a_3 = -1$ . Suppose also that for  $n \geq 4$  it is known that  $a_n = a_{n-1} + 2a_{n-2} + 3a_{n-3}$ . Find  $a_7$ . Are you able to guess a formula for the  $k^{\text{th}}$  term of this sequence?
- Let  $f: \{t \in \mathbb{R} : t \neq -1\} \rightarrow \mathbb{R}$  be defined by  $f(t) \equiv \frac{t}{t+1}$ . Find  $f^{-1}$  if possible.
- A function,  $f: \mathbb{R} \rightarrow \mathbb{R}$  is a strictly increasing function if whenever  $x < y$ , it follows that  $f(x) < f(y)$ . If  $f$  is a strictly increasing function, does  $f^{-1}$  always exist? Explain your answer.

7. Let  $f(t)$  be defined by

$$f(t) = \begin{cases} 2t + 1 & \text{if } t \leq 1 \\ t & \text{if } t > 1 \end{cases}.$$

Find  $f^{-1}$  if possible.

8. Suppose  $f : D(f) \rightarrow R(f)$  is one to one,  $R(f) \subseteq D(g)$ , and  $g : D(g) \rightarrow R(g)$  is one to one. Does it follow that  $g \circ f$  is one to one?

9. If  $f : \mathbb{R} \rightarrow \mathbb{R}$  and  $g : \mathbb{R} \rightarrow \mathbb{R}$  are two one to one functions, which of the following are necessarily one to one on their domains? Explain why or why not by giving a proof or an example.

(a)  $f + g$

(b)  $fg$

(c)  $f^3$

(d)  $f/g$

10. Draw the graph of the function  $f(x) = x^3 + 1$ .

11. Draw the graph of the function  $f(x) = x^2 + 2x + 2$ .

12. Draw the graph of the function  $f(x) = \frac{x}{1+x}$ .

13. Suppose  $a_n = \frac{1}{n}$  and let  $n_k = 2^k$ . Find  $b_k$  where  $b_k = a_{n_k}$ .

14. If  $X_i$  are sets and for some  $j$ ,  $X_j = \emptyset$ , the empty set. Verify carefully that  $\prod_{i=1}^n X_i = \emptyset$ .

15. Suppose  $f(x) + f\left(\frac{1}{x}\right) = 7x$  and  $f$  is a function defined on  $\mathbb{R} \setminus \{0\}$ , the nonzero real numbers. Find all values of  $x$  where  $f(x) = 1$  if there are any. Does there exist any such function?

16. Does there exist a function  $f$ , satisfying  $f(x) - f\left(\frac{1}{x}\right) = 3x$  which has both  $x$  and  $\frac{1}{x}$  in the domain of  $f$ ?

17. In the situation of the Fibonacci sequence show that the formula for the  $n^{\text{th}}$  term can be found and is given by

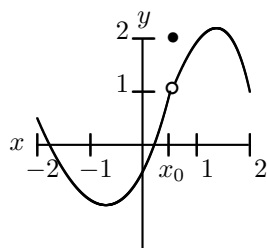
$$a_n = \frac{\sqrt{5}}{5} \left( \frac{1 + \sqrt{5}}{2} \right)^n - \frac{\sqrt{5}}{5} \left( \frac{1 - \sqrt{5}}{2} \right)^n.$$

**Hint:** You might be able to do this by induction but a better way would be to look for a solution to the recurrence relation,  $a_{n+2} \equiv a_n + a_{n+1}$  of the form  $r^n$ . You will be able to show that there are two values of  $r$  which work, one of which is  $r = \frac{1+\sqrt{5}}{2}$ . Next you can observe that if  $r_1^n$  and  $r_2^n$  both satisfy the recurrence relation then so does  $cr_1^n + dr_2^n$  for any choice of constants  $c, d$ . Then you try to pick  $c$  and  $d$  such that the conditions,  $a_1 = 1$  and  $a_2 = 1$  both hold.

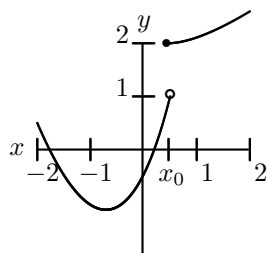
## 5.3 Continuous Functions

The concept of function is far too general to be useful by itself. There are various ways to restrict the concept in order to study something interesting and the types of restrictions considered depend very much on what you find interesting. In Calculus, the most fundamental restriction made is to assume the functions are continuous. Continuous functions are those in which a sufficiently small change in  $x$  results in a small change in  $f(x)$ . They rule out things which could never happen physically. For example, it is not possible for a car to jump from one point to another instantly. Making this restriction precise turns out to be surprisingly difficult although many of the most important theorems about continuous functions seem intuitively clear.

Before giving the careful mathematical definitions, here are examples of graphs of functions which are not continuous at the point  $x_0$ .



You see, there is a hole in the picture of the graph of this function and instead of filling in the hole with the appropriate value,  $f(x_0)$  is too large. This is called a removable discontinuity because the problem can be fixed by redefining the function at the point  $x_0$ . Here is another example.



You see from this picture that there is no way to get rid of the jump in the graph of this function by simply redefining the value of the function at  $x_0$ . That is why it is called a nonremovable discontinuity or jump discontinuity. Now that pictures have been given of what it is desired to eliminate, it is time to give the precise definition.

The definition which follows, due to Cauchy<sup>1</sup> and Weierstrass<sup>2</sup> is the precise way to exclude the sort of behavior described above and all statements about continuous functions must ultimately rest on this definition from now on.

**Definition 5.3.1** *A function  $f : D(f) \subseteq \mathbb{R} \rightarrow \mathbb{R}$  is continuous at  $x \in D(f)$  if for each  $\varepsilon > 0$  there exists  $\delta > 0$  such that whenever  $y \in D(f)$  and*

$$|y - x| < \delta$$

*it follows that*

$$|f(x) - f(y)| < \varepsilon.$$

*A function,  $f$  is continuous if it is continuous at every point of  $D(f)$ .*

In sloppy English this definition says roughly the following: A function,  $f$  is continuous at  $x$  when it is possible to make  $f(y)$  as close as desired to  $f(x)$  provided  $y$  is taken close enough to  $x$ . In fact this statement in words is pretty much the way Cauchy described it. The completely rigorous definition above is due to Weierstrass. This definition does indeed rule out the sorts of graphs drawn above. Consider the second nonremovable discontinuity.

---

<sup>1</sup>Augustin Louis Cauchy 1789-1857 was the son of a lawyer who was married to an aristocrat. He was born in France just after the fall of the Bastille and his family fled the reign of terror and hid in the countryside till it was over. Cauchy was educated at first by his father who taught him Greek and Latin. Eventually Cauchy learned many languages. He was also a good Catholic.

After the reign of terror, the family returned to Paris and Cauchy studied at the university to be an engineer but became a mathematician although he made fundamental contributions to physics and engineering. Cauchy was one of the most prolific mathematicians who ever lived. He wrote several hundred papers which fill 24 volumes. He also did research on many topics in mechanics and physics including elasticity, optics and astronomy. More than anyone else, Cauchy invented the subject of complex analysis. He is also credited with giving the first rigorous definition of continuity.

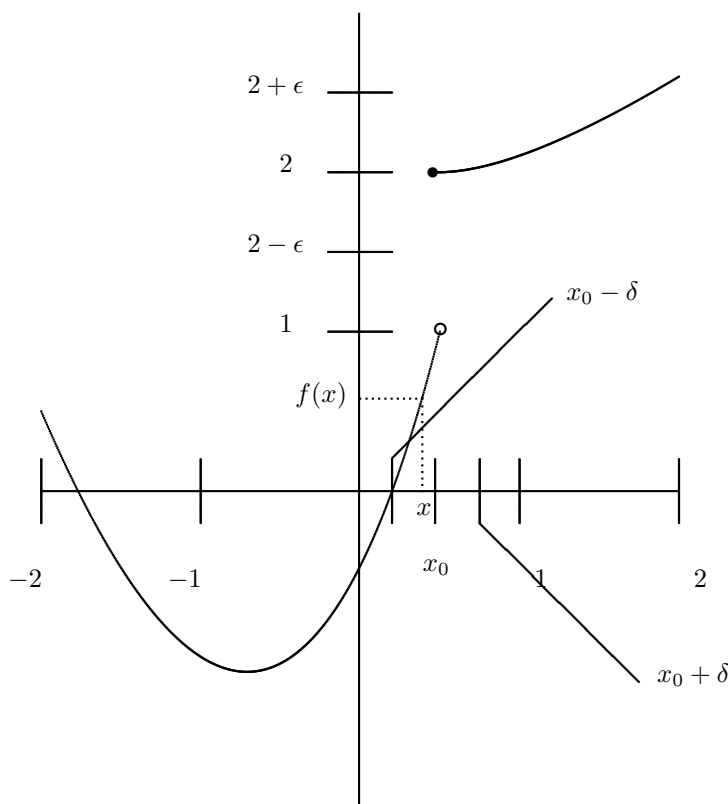
He married in 1818 and lived for 12 years with his wife and two daughters in Paris till the revolution of 1830. Cauchy refused to take the oath of allegiance to the new ruler and ended up leaving his family and going into exile for 8 years.

Notwithstanding his great achievements he was not known as a popular teacher.

<sup>2</sup>Wilhelm Theodor Weierstrass 1815-1897 brought calculus to essentially the state it is in now. When he was a secondary school teacher, he wrote a paper which was so profound that he was granted a doctor's degree. He made fundamental contributions to partial differential equations, complex analysis, calculus of variations, and many other topics. He also discovered some pathological examples such as space filling curves. Cauchy gave the definition in words and Weierstrass, somewhat later produced the totally rigorous  $\varepsilon \delta$  definition presented here. The need for rigor in the subject of calculus was only realized over a long period of time.



The removable discontinuity case is similar.



For the  $\epsilon$  shown you can see from the picture that no matter how small you take  $\delta$ , there will be points,  $x$ , between  $x_0 - \delta$  and  $x_0$  where  $f(x) < 2 + \epsilon$ . In particular, for these values of  $x$ ,  $|f(x) - f(x_0)| > \epsilon$ . Therefore, the definition of continuity given above excludes the situation in which there is a jump in the function. Similar reasoning shows it excludes the removable discontinuity case as well. There are many ways a function can fail to be continuous and it is impossible to list them all by drawing pictures. This is why it is so important to use the definition. The other thing to notice is that the concept of continuity as described in the definition is a point property. That is to say it is a property which a function may or may not have at a single point. Here is an example.

**Example 5.3.2** Let

$$f(x) = \begin{cases} x & \text{if } x \text{ is rational} \\ 0 & \text{if } x \text{ is irrational} \end{cases}.$$

*This function is continuous at  $x = 0$  and nowhere else.*

To verify the assertion about the above function, first show it is not continuous at  $x$  if  $x \neq 0$ . Take such an  $x$  and let  $\epsilon = |x|/2$ . Now let  $\delta > 0$  be completely arbitrary. In the interval,  $(x - \delta, x + \delta)$  there are rational numbers,  $y_1$  such that  $|y_1| > |x|$  and irrational numbers,  $y_2$ . Thus  $|f(y_1) - f(y_2)| = |y_1| > |x|$ . If  $f$  were continuous at  $x$ , there would exist  $\delta > 0$  such that for every point,  $y \in (x - \delta, x + \delta)$ ,  $|f(y) - f(x)| < \epsilon$ . But then, letting  $y_1$

and  $y_2$  be as just described,

$$\begin{aligned} |x| &< |y_1| = |f(y_1) - f(y_2)| \\ &\leq |f(y_1) - f(x)| + |f(x) - f(y_2)| < 2\varepsilon = |x|, \end{aligned}$$

which is a contradiction. Since a contradiction is obtained by assuming that  $f$  is continuous at  $x$ , it must be concluded that  $f$  is not continuous there. To see  $f$  is continuous at 0, let  $\varepsilon > 0$  be given and let  $\delta = \varepsilon$ . Then if  $|y - 0| < \delta = \varepsilon$ , Then

$$\begin{aligned} |f(y) - f(0)| &= 0 \text{ if } y \text{ is irrational} \\ |f(y) - f(0)| &= |y| < \varepsilon \text{ if } y \text{ is rational.} \end{aligned}$$

either way, whenever  $|y - 0| < \delta$ , it follows  $|f(y) - f(0)| < \varepsilon$  and so  $f$  is continuous at  $x = 0$ . How did I know to let  $\delta = \varepsilon$ ? That is a very good question. The choice of  $\delta$  for a particular  $\varepsilon$  is usually arrived at by using intuition, the actual  $\varepsilon \delta$  argument reduces to a verification that the intuition was correct. Here is another example.

**Example 5.3.3** Show the function,  $f(x) = -5x + 10$  is continuous at  $x = -3$ .

To do this, note first that  $f(-3) = 25$  and it is desired to verify the conditions for continuity. Consider the following.

$$|-5x + 10 - (25)| = 5|x - (-3)|.$$

This allows one to find a suitable  $\delta$ . If  $\varepsilon > 0$  is given, let  $0 < \delta \leq \frac{1}{5}\varepsilon$ . Then if  $0 < |x - (-3)| < \delta$ , it follows from this inequality that

$$|-5x + 10 - (25)| = 5|x - (-3)| < 5\frac{1}{5}\varepsilon = \varepsilon.$$

Sometimes the determination of  $\delta$  in the verification of continuity can be a little more involved. Here is another example.

**Example 5.3.4** Show the function,  $f(x) = \sqrt{2x + 12}$  is continuous at  $x = 5$ .

First note  $f(5) = \sqrt{22}$ . Now consider:

$$\begin{aligned} \left| \sqrt{2x + 12} - \sqrt{22} \right| &= \left| \frac{2x + 12 - 22}{\sqrt{2x + 12} + \sqrt{22}} \right| \\ &= \frac{2}{\sqrt{2x + 12} + \sqrt{22}} |x - 5| \leq \frac{1}{11} \sqrt{22} |x - 5| \end{aligned}$$

whenever  $|x - 5| < 1$  because for such  $x$ ,  $\sqrt{2x + 12} > 0$ . Now let  $\varepsilon > 0$  be given. Choose  $\delta$  such that  $0 < \delta \leq \min\left(1, \frac{\varepsilon\sqrt{22}}{2}\right)$ . Then if  $|x - 5| < \delta$ , all the inequalities above hold and

$$\left| \sqrt{2x + 12} - \sqrt{22} \right| \leq \frac{2}{\sqrt{22}} |x - 5| < \frac{2}{\sqrt{22}} \frac{\varepsilon\sqrt{22}}{2} = \varepsilon.$$

**Exercise 5.3.5** Show  $f(x) = -3x^2 + 7$  is continuous at  $x = 7$ .

First observe  $f(7) = -140$ . Now

$$|-3x^2 + 7 - (-140)| = 3|x + 7||x - 7| \leq 3(|x| + 7)|x - 7|.$$

If  $|x - 7| < 1$ , it follows from the version of the triangle inequality which states  $||s| - |t|| \leq |s - t|$  that  $|x| < 1 + 7$ . Therefore, if  $|x - 7| < 1$ , it follows that

$$\begin{aligned} |-3x^2 + 7 - (-140)| &\leq 3((1 + 7) + 7)|x - 7| \\ &= 3(1 + 27)|x - 7| = 84|x - 7|. \end{aligned}$$

Now let  $\varepsilon > 0$  be given. Choose  $\delta$  such that  $0 < \delta \leq \min(1, \frac{\varepsilon}{84})$ . Then for  $|x - 7| < \delta$ , it follows

$$|-3x^2 + 7 - (-140)| \leq 84|x - 7| < 84\left(\frac{\varepsilon}{84}\right) = \varepsilon.$$

These  $\varepsilon \delta$  proofs will not be emphasized any more than necessary. However, you should try a few of them because until you master this concept, you will not really understand calculus as it has been understood for approximately 150 years. The best you can do without this definition is to gain an understanding of the subject as it was understood by people in the 1700's, before the need for rigor was realized. Don't be discouraged by these historical observations. If you are able to master calculus as understood by Lagrange or Laplace<sup>3</sup>, you will have learned some very profound ideas even if they did originate in the eighteenth century.

## 5.4 Sufficient Conditions For Continuity

The next theorem is a fundamental result which will allow us to worry less about the  $\varepsilon \delta$  definition of continuity.

**Theorem 5.4.1** *The following assertions are valid for  $f, g$  functions and  $a, b$  numbers.*

1. *The function,  $af + bg$  is continuous at  $x$  when  $f, g$  are continuous at  $x \in D(f) \cap D(g)$  and  $a, b \in \mathbb{R}$ .*
2. *If  $f$  and  $g$  are each real valued functions continuous at  $x$ , then  $fg$  is continuous at  $x$ . If, in addition to this,  $g(x) \neq 0$ , then  $f/g$  is continuous at  $x$ .*
3. *If  $f$  is continuous at  $x$ ,  $f(x) \in D(g) \subseteq \mathbb{R}$ , and  $g$  is continuous at  $f(x)$ , then  $g \circ f$  is continuous at  $x$ .*
4. *The function  $f: \mathbb{R} \rightarrow \mathbb{R}$ , given by  $f(x) = |x|$  is continuous.*

The proof of this theorem is in the last section of this chapter but its conclusions are not surprising. For example the first claim says that  $(af + bg)(y)$  is close to  $(af + bg)(x)$  when  $y$  is close to  $x$  provided the same can be said about  $f$  and  $g$ . For the third claim, continuity of  $f$  indicates that if  $y$  is close enough to  $x$  then  $f(x)$  is close to  $f(y)$  and so by continuity of  $g$  at  $f(x)$ ,  $g(f(y))$  is close to  $g(f(x))$ . The fourth claim is verified as follows.

$$|x| = |x - y + y| \leq |x - y| + |y|$$

and so

$$|x| - |y| \leq |x - y|.$$

---

<sup>3</sup>Lagrange and Laplace were two great physicists of the 1700's. They made fundamental contributions to the calculus of variations and to mechanics and astronomy.

Similarly,

$$|y| - |x| \leq |x - y|.$$

Therefore,

$$||x| - |y|| \leq |x - y|.$$

It follows that if  $\varepsilon > 0$  is given, one can take  $\delta = \varepsilon$  and obtain that for  $|x - y| < \delta = \varepsilon$ ,

$$||x| - |y|| \leq |x - y| < \delta = \varepsilon$$

which shows continuity of the function,  $f(x) = |x|$ .

## 5.5 Continuity Of Circular Functions

The functions  $\sin x$  and  $\cos x$  are often called the circular functions. This is because for each  $x \in \mathbb{R}$ ,  $(\cos x, \sin x)$  is a point on the unit circle.

**Theorem 5.5.1** *The functions,  $\cos$  and  $\sin$  are continuous.*

**Proof:** First it will be shown that  $\cos$  and  $\sin$  are continuous at 0. By Corollary 3.5.4 on Page 59 the following inequality is valid for small positive values of  $\theta$ .

$$1 - \cos \theta + \sin \theta \geq \theta \geq \sin \theta.$$

It follows that for  $\theta$  small and positive,  $|\theta| \geq |\sin \theta| = \sin \theta$ . If  $\theta < 0$ , then  $-\theta = |\theta| > 0$  and  $-\theta \geq \sin(-\theta)$ . But then this means  $|\sin \theta| = -\sin \theta = \sin(-\theta) \leq -\theta = |\theta|$  in this case also. Therefore, whenever  $|\theta|$  is small enough,

$$|\theta| \geq |\sin \theta|.$$

Now let  $\varepsilon > 0$  be given and take  $\delta = \varepsilon$ . Then if  $|\theta| < \delta$ , it follows

$$|\sin \theta - 0| = |\sin \theta - \sin 0| = |\sin \theta| \leq |\theta| < \delta = \varepsilon,$$

showing  $\sin$  is continuous at 0.

Next, note that for  $|\theta| < \pi/2$ ,  $\cos \theta \geq 0$  and so for such  $\theta$ ,

$$\sin^2 \theta \geq \frac{\sin^2 \theta}{1 + \cos \theta} = \frac{1 - \cos^2 \theta}{1 + \cos \theta} = 1 - \cos \theta \geq 0. \quad (5.2)$$

From the first part of this argument for  $\sin$ , given  $\varepsilon > 0$  there exists  $\delta > 0$  such that if  $|\theta| < \delta$ , then  $|\sin \theta| < \sqrt{\varepsilon}$ . It follows from (5.2) that if  $|\theta| < \delta$ , then  $\varepsilon > 1 - \cos \theta \geq 0$ . This proves these functions are continuous at 0. Now  $y = (y - x) + x$  and so

$$\cos y = \cos(y - x) \cos x - \sin(y - x) \sin x.$$

Therefore,

$$\cos y - \cos x = \cos(y - x) \cos x - \sin(y - x) \sin x - \cos x$$

and so, since  $|\cos x|, |\sin x| \leq 1$ ,

$$\begin{aligned} |\cos y - \cos x| &\leq |\cos x (\cos(y - x) - 1)| + |\sin x| |\sin(y - x)| \\ &\leq |\cos(y - x) - 1| + |\sin(y - x)|. \end{aligned}$$

From the first part of this theorem, if  $|y - x|$  is sufficiently small, both of these last two terms are less than  $\varepsilon/2$  and this proves  $\cos$  is continuous at  $x$ . The proof that  $\sin$  is continuous is left for you to verify.

## 5.6 Exercises

1. Let  $f(x) = 2x + 7$ . Show  $f$  is continuous at every point  $x$ . **Hint:** You need to let  $\varepsilon > 0$  be given. In this case, you should try  $\delta \leq \varepsilon/2$ . Note that if one  $\delta$  works in the definition, then so does any smaller  $\delta$ .

2. Let  $f(x) = x^2 + 1$ . Show  $f$  is continuous at  $x = 3$ . **Hint:**

$$\begin{aligned} |f(x) - f(3)| &= |x^2 + 1 - (9 + 1)| \\ &= |x + 3| |x - 3|. \end{aligned}$$

Thus if  $|x - 3| < 1$ , it follows from the triangle inequality,  $|x| < 1 + 3 = 4$  and so

$$|f(x) - f(3)| < 4|x - 3|.$$

Now try to complete the argument by letting  $\delta \leq \min(1, \varepsilon/4)$ . The symbol,  $\min$  means to take the minimum of the two numbers in the parenthesis.

3. Let  $f(x) = x^2 + 1$ . Show  $f$  is continuous at  $x = 4$ .
4. Let  $f(x) = 2x^2 + 1$ . Show  $f$  is continuous at  $x = 1$ .
5. Let  $f(x) = x^2 + 2x$ . Show  $f$  is continuous at  $x = 2$ . Then show it is continuous at every point.
6. Let  $f(x) = |2x + 3|$ . Show  $f$  is continuous at every point. **Hint:** Review the two versions of the triangle inequality for absolute values.
7. Let  $f(x) = \frac{1}{x^2 + 1}$ . Show  $f$  is continuous at every value of  $x$ .
8. Show  $\sin$  is continuous.
9. Let  $f(x) = \sqrt{x}$  show  $f$  is continuous at every value of  $x$  in its domain. **Hint:** You might want to make use of the identity,

$$\sqrt{x} - \sqrt{y} = \frac{x - y}{\sqrt{x} + \sqrt{y}}$$

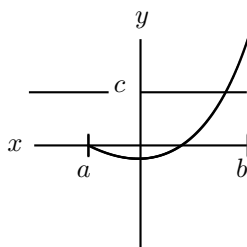
at some point in your argument.

10. Using Theorem 5.4.1, show all polynomials are continuous and that a rational function is continuous at every point of its domain. **Hint:** First show the function given as  $f(x) = x$  is continuous and then use the theorem.
11. Let  $f(x) = \begin{cases} 1 & \text{if } x \in \mathbb{Q} \\ 0 & \text{if } x \notin \mathbb{Q} \end{cases}$  and consider  $g(x) = f(x)\sin x$ . Determine where  $g$  is continuous and explain your answer.

## 5.7 Properties Of Continuous Functions

Continuous functions have many important properties which are consequences of the completeness axiom. Proofs of these theorems are in the last section at the end of this chapter.

The next theorem is called the intermediate value theorem and the following picture illustrates its conclusion. It gives the existence of a certain point.



You see in the picture there is a horizontal line,  $y = c$  and a continuous function which starts off less than  $c$  at the point  $a$  and ends up greater than  $c$  at point  $b$ . The intermediate value theorem says there is some point between  $a$  and  $b$  such that the value of the function at this point equals  $c$ . You see this taking place in the above picture where the line and the graph of the function cross. The  $x$  value at this point is the one whose existence is guaranteed by the theorem. It may seem this is obvious but without completeness the conclusion of the theorem cannot be drawn. Nevertheless, the above picture makes this theorem very easy to believe.

**Theorem 5.7.1** Suppose  $f : [a, b] \rightarrow \mathbb{R}$  is continuous and suppose  $f(a) < c < f(b)$ . Then there exists  $x \in (a, b)$  such that  $f(x) = c$ .

**Example 5.7.2** Does there exist a solution to the equation  $\sqrt{x^4 + 7} - x^3 \sin x = 0$ ?

By Theorem 5.4.1 and Problem 9 on Page 93 it follows easily that the function,  $f$ , given by  $f(x) = \sqrt{x^4 + 7} - x^3 \sin x$  is continuous. Also,  $f(0) = \sqrt{7} > 0$  while

$$f\left(\frac{5\pi}{2}\right) = \sqrt{\left(\frac{5\pi}{2}\right)^4 + 7} - \left(\frac{5\pi}{2}\right)^3 \sin\left(\frac{5\pi}{2}\right)$$

which is approximately equal to  $-422.7313318316316 < 0$ . Therefore, by the intermediate value theorem there must exist  $x \in (0, \frac{5\pi}{2})$  such that  $f(x) = 0$ .

This example illustrates the use of this major theorem very well. It says something exists but it does not tell how to find it.

**Definition 5.7.3** A function,  $f$ , defined on some interval is strictly increasing if whenever  $x < y$ , it follows  $f(x) < f(y)$ . The function is strictly decreasing if whenever  $x < y$ , it follows  $f(x) > f(y)$ .

You should draw a picture of the graph of a strictly increasing or decreasing function from the definition.

**Lemma 5.7.4** Let  $\phi : [a, b] \rightarrow \mathbb{R}$  be a one to one continuous function. Then  $\phi$  is either strictly increasing or strictly decreasing.

This lemma is not real easy to prove but it is one of those things which seems obvious. To say a function is one to one is to say that every horizontal line intersects the graph of the function in no more than one point. (This is called the horizontal line test.) Now if your function is continuous (having no jumps) and is one to one, try to imagine how this could happen without it being either strictly increasing or decreasing and you will soon see this is highly believable and in fact, for it to fail would be incredible. The proof of this lemma is in the last section of this chapter in case you are interested.

**Corollary 5.7.5** *Let  $\phi : (a, b) \rightarrow \mathbb{R}$  be a one to one continuous function. Then  $\phi$  is either strictly increasing or strictly decreasing.*

The proof of this corollary is the same as the proof of the lemma. The next corollary follows from the above.

**Corollary 5.7.6** *Let  $f : (a, b) \rightarrow \mathbb{R}$  be one to one and continuous. Then  $f(a, b)$  is an open interval,  $(c, d)$  and  $f^{-1} : (c, d) \rightarrow (a, b)$  is continuous. Also, if  $f : [a, b] \rightarrow \mathbb{R}$  is one to one and continuous, then  $f([a, b])$  is a closed interval,  $[c, d]$  and  $f^{-1} : [c, d] \rightarrow [a, b]$  is continuous.*

This corollary is not too surprising either. To view the graph of the inverse function, simply turn things on the side and switch  $x$  and  $y$ . If the original graph has no jumps in it, neither will the new graph. Of course, the concept of continuity is tied to a rigorous definition, not to the drawing of pictures. This is why there is a proof in the last section of this chapter.

In Russia there is a kind of doll called a matrushka doll. You pick it up and notice it comes apart in the center. Separating the two halves you find an identical doll inside. Then you notice this inside doll also comes apart in the center. Separating the two halves, you find yet another identical doll inside. This goes on quite a while until the final doll is in one piece. The nested interval lemma is like a matrushka doll except the process never stops. It involves a sequence of intervals, the first containing the second, the second containing the third, the third containing the fourth and so on. The fundamental question is whether there exists a point in all the intervals.

**Lemma 5.7.7** *Let  $I_k = [a^k, b^k]$  and suppose that for all  $k = 1, 2, \dots$ ,*

$$I_k \supseteq I_{k+1}.$$

*Then there exists a point,  $c \in \mathbb{R}$  which is an element of every  $I_k$ .*

**Proof:** Since  $I_k \supseteq I_{k+1}$ , this implies

$$a^k \leq a^{k+1}, \quad b^k \geq b^{k+1}. \quad (5.3)$$

Consequently, if  $k \leq l$ ,

$$a^l \leq a^k \leq b^k \leq b^l. \quad (5.4)$$

Now define

$$c \equiv \sup \{a^l : l = 1, 2, \dots\}$$

By the first inequality in (5.3), and (5.4)

$$a^k \leq c = \sup \{a^l : l = k, k+1, \dots\} \leq b^k \quad (5.5)$$

for each  $k = 1, 2, \dots$ . Thus  $c \in I_k$  for every  $k$  and this proves the lemma. If this went too fast, the reason for the last inequality in (5.5) is that from (5.4),  $b^k$  is an upper bound to  $\{a^l : l = k, k+1, \dots\}$ . Therefore, it is at least as large as the least upper bound.

This is really quite a remarkable result and may not seem so obvious. Consider the intervals  $I_k \equiv (0, 1/k)$ . Then there is no point which lies in all these intervals because no negative number can be in all the intervals and  $1/k$  is smaller than a given positive number whenever  $k$  is large enough. Thus the only candidate for being in all the intervals is 0 and 0 has been left out of them all. The problem here is that the endpoints of the intervals were not included contrary to the hypotheses of the above lemma in which all the intervals included the endpoints.

With the nested interval lemma, it becomes possible to prove the following lemma which shows a function continuous on a closed interval in  $\mathbb{R}$  is bounded.

**Lemma 5.7.8** *Let  $I = [a, b]$  and let  $f : I \rightarrow \mathbb{R}$  be continuous. Then  $f$  is bounded. That is there exist numbers,  $m$  and  $M$  such that for all  $x \in [a, b]$ ,*

$$m \leq f(x) \leq M.$$

**Proof:** Let  $I \equiv I_0$  and suppose  $f$  is not bounded on  $I_0$ . Consider the two sets,  $[a, \frac{a+b}{2}]$  and  $[\frac{a+b}{2}, b]$ . Since  $f$  is not bounded on  $I_0$ , it follows that  $f$  must fail to be bounded on at least one of these sets. Let  $I_1$  be one of these on which  $f$  is not bounded. Now do to  $I_1$  what was done to  $I_0$  to obtain  $I_2 \subseteq I_1$  and for any two points,  $x, y \in I_2$

$$|x - y| \leq 2^{-1} \frac{b-a}{2} \leq 2^{-2} (b-a).$$

Continue in this way obtaining sets,  $I_k$  such that  $I_k \supseteq I_{k+1}$  and for any two points in  $I_k$ ,  $x, y$ ,  $|x - y| \leq 2^{-k} (b-a)$ . By the nested interval lemma, there exists a point,  $c$  which is contained in each  $I_k$ . Also, by continuity, there exists a  $\delta > 0$  such that if  $|c - y| < \delta$ , then

$$|f(c) - f(y)| < 1. \quad (5.6)$$

Let  $k$  be so large that  $2^{-k} (b-a) < \delta$ . Then for every  $y \in I_k$ ,  $|c - y| < \delta$  and so (5.6) holds for all such  $y$ . But this implies that for all  $y \in I_k$ ,

$$|f(y)| \leq |f(c)| + 1$$

which shows that  $f$  is bounded on  $I_k$  contrary to the way  $I_k$  was chosen. This contradiction proves the lemma.

**Example 5.7.9** *Let  $f(x) = 1/x$  for  $x \in (0, 1)$ .*

Clearly,  $f$  is not bounded. Does this violate the conclusion of the nested interval lemma? It does not because the end points of the interval involved are not in the interval. The same function defined on  $[.000001, 1)$  would have been bounded although in this case the boundedness of the function would not follow from the above lemma because it fails to include the right endpoint.

The next theorem is known as the max min theorem or extreme value theorem.

**Theorem 5.7.10** *Let  $I = [a, b]$  and let  $f : I \rightarrow \mathbb{R}$  be continuous. Then  $f$  achieves its maximum and its minimum on  $I$ . This means there exist,  $x_1, x_2 \in I$  such that for all  $x \in I$ ,*

$$f(x_1) \leq f(x) \leq f(x_2).$$

**Proof:** By completeness of  $\mathbb{R}$  and Lemma 5.7.8  $f(I)$  has a least upper bound,  $M$ . If for all  $x \in I$ ,  $f(x) \neq M$ , then by Theorem 5.4.1, the function,  $g(x) \equiv (M - f(x))^{-1} = \frac{1}{M - f(x)}$  is continuous on  $I$ . Since  $M$  is the least upper bound of  $f(I)$  there exist points,  $x \in I$  such that  $(M - f(x))$  is as small as desired. Consequently,  $g$  is not bounded above, contrary to Lemma 5.7.8. Therefore, there must exist some  $x \in I$  such that  $f(x) = M$ . This proves  $f$  achieves its maximum. The argument for the minimum is similar. Alternatively, you could consider the function  $h(x) = M - f(x)$ . Then use what was just proved to conclude  $h$  achieves its maximum at some point,  $x_1$ . Thus  $h(x_1) \geq h(x)$  for all  $x \in I$  and so  $M - f(x_1) \geq M - f(x)$  for all  $x \in I$  which implies  $f(x_1) \leq f(x)$  for all  $x \in I$ . This proves the theorem.



## 5.8 Exercises

1. Give an example of a continuous function defined on  $(0, 1)$  which does not achieve its maximum on  $(0, 1)$ .
2. Give an example of a continuous function defined on  $(0, 1)$  which is bounded but which does not achieve either its maximum or its minimum.
3. Give an example of a discontinuous function defined on  $[0, 1]$  which is bounded but does not achieve either its maximum or its minimum.
4. Give an example of a continuous function defined on  $[0, 1) \cup (1, 2]$  which is positive at 2, negative at 0 but is not equal to zero for any value of  $x$ .
5. Give an example of a function which is one to one but neither strictly increasing nor strictly decreasing. **Hint:** Look for discontinuous functions satisfying the horizontal line test.
6. Do you believe in  $\sqrt[7]{8}$ ? That is, does there exist a number which multiplied by itself seven times yields 8? Before you jump to any conclusions, the number you get on your calculator is wrong. In fact, your calculator does not even know about  $\sqrt[7]{8}$ . All it can do is try to approximate it and what it gives you is this approximation. Why does it exist? **Hint:** Use the intermediate value theorem on the function,  $f(x) = x^7 - 8$ .
7. Let  $f(x) = x - \sqrt{2}$  for  $x \in \mathbb{Q}$ , the rational numbers. Show that even though  $f(0) < 0$  and  $f(2) > 0$ , there is no point in  $\mathbb{Q}$  where  $f(x) = 0$ . Does this contradict the intermediate value theorem? Explain.
8. It has been known since the time of Pythagoras that  $\sqrt{2}$  is irrational. If you throw out all the irrational numbers, show that the conclusion of the intermediate value theorem could no longer be obtained. That is, show there exists a function which starts off less than zero and ends up larger than zero and yet there is no number where the function equals zero. **Hint:** Try  $f(x) = x^2 - 2$ . You supply the details.

## 5.9 Limits Of A Function

A concept closely related to continuity is that of the limit of a function.

**Definition 5.9.1** Let  $f : D(f) \subseteq \mathbb{R} \rightarrow \mathbb{R}$  be a function where  $D(f) \supseteq (x - r, x) \cup (x, x + r)$  for some  $r > 0$ . Note that  $f$  is not necessarily defined at  $x$ . Then

$$\lim_{y \rightarrow x} f(y) = L$$

if and only if the following condition holds. For all  $\varepsilon > 0$  there exists  $\delta > 0$  such that if

$$0 < |y - x| < \delta,$$

then,

$$|L - f(y)| < \varepsilon.$$

If everything is the same as the above, except  $y$  is required to be larger than  $x$  and  $f$  is only required to be defined on  $(x, x + r)$ , then the notation is

$$\lim_{y \rightarrow x^+} f(y) = L.$$

If  $f$  is only required to be defined on  $(x - r, x)$  and  $y$  is required to be less than  $x$ , with the same conditions above, we write

$$\lim_{y \rightarrow x^-} f(y) = L.$$

Limits are also taken as a variable “approaches” infinity. Of course nothing is “close” to infinity and so this requires a slightly different definition.

$$\lim_{x \rightarrow \infty} f(x) = L$$

if for every  $\varepsilon > 0$  there exists  $l$  such that whenever  $x > l$ ,

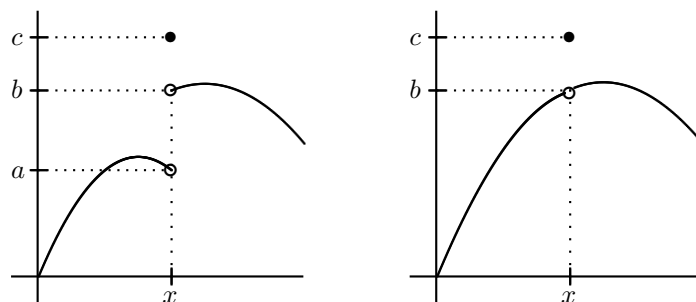
$$|f(x) - L| < \varepsilon \quad (5.7)$$

and

$$\lim_{x \rightarrow -\infty} f(x) = L$$

if for every  $\varepsilon > 0$  there exists  $l$  such that whenever  $x < l$ , (5.7) holds.

The following pictures illustrate some of these definitions.



In the left picture is shown the graph of a function. Note the value of the function at  $x$  equals  $c$  while  $\lim_{y \rightarrow x+} f(y) = b$  and  $\lim_{y \rightarrow x-} f(y) = a$ . In the second picture,  $\lim_{y \rightarrow x} f(y) = b$ . Note that the value of the function at the point  $x$  has nothing to do with the limit of the function in any of these cases. **The value of a function at  $x$  is irrelevant to the value of the limit at  $x$ !** This must always be kept in mind. You do not evaluate interesting limits by computing  $f(x)$ ! In the above picture,  $f(x)$  is always wrong! It may be the case that  $f(x)$  is right but this is merely a happy coincidence when it occurs and as explained below in Theorem 5.9.6, this is equivalent to  $f$  being continuous at  $x$ .

**Theorem 5.9.2** If  $\lim_{y \rightarrow x} f(y) = L$  and  $\lim_{y \rightarrow x} f(y) = L_1$ , then  $L = L_1$ .

**Proof:** Let  $\varepsilon > 0$  be given. There exists  $\delta > 0$  such that if  $0 < |y - x| < \delta$ , then

$$|f(y) - L| < \varepsilon, \quad |f(y) - L_1| < \varepsilon.$$

Therefore, for such  $y$ ,

$$|L - L_1| \leq |L - f(y)| + |f(y) - L_1| < \varepsilon + \varepsilon = 2\varepsilon.$$

Since  $\varepsilon > 0$  was arbitrary, this shows  $L = L_1$ .

The above theorem holds for any of the kinds of limits presented in the above definition.

Another concept is that of a function having either  $\infty$  or  $-\infty$  as a limit. In this case, the values of the function do not ever get close to their target because nothing can be close to  $\pm\infty$ . Roughly speaking, the limit of the function equals  $\infty$  if the values of the function are ultimately larger than any given number. More precisely:

**Definition 5.9.3** If  $f(x) \in \mathbb{R}$ , then  $\lim_{y \rightarrow x} f(x) = \infty$  if for every number  $l$ , there exists  $\delta > 0$  such that whenever  $|y - x| < \delta$ , then  $f(x) > l$ .  $\lim_{x \rightarrow \infty} f(x) = \infty$  if for all  $k$ , there exists  $l$  such that  $f(x) > k$  whenever  $x > l$ . One sided limits and limits as the variable approaches  $-\infty$ , are defined similarly.

It may seem there is a lot to memorize here. In fact, this is not so because all the definitions are intuitive when you understand them.

**Theorem 5.9.4** In this theorem, the symbol,  $\lim_{y \rightarrow x}$  denotes any of the limits described above. Suppose  $\lim_{y \rightarrow x} f(y) = L$  and  $\lim_{y \rightarrow x} g(y) = K$  where  $K$  and  $L$  are real numbers in  $\mathbb{R}$ . Then if  $a, b \in \mathbb{R}$ ,

$$\lim_{y \rightarrow x} (af(y) + bg(y)) = aL + bK, \quad (5.8)$$

$$\lim_{y \rightarrow x} fg(y) = LK \quad (5.9)$$

and if  $K \neq 0$ ,

$$\lim_{y \rightarrow x} \frac{f(y)}{g(y)} = \frac{L}{K}. \quad (5.10)$$

Also, if  $h$  is a continuous function defined near  $L$ , then

$$\lim_{y \rightarrow x} h \circ f(y) = h(L). \quad (5.11)$$

Suppose  $\lim_{y \rightarrow x} f(y) = L$ . If  $f(y) \leq a$  all  $y$  of interest, then  $L \leq a$  and if  $f(y) \geq a$  then  $L \geq a$ .

**Proof:** The proof of (5.8) is left for you. It is like a corresponding theorem for continuous functions. Next consider (5.9). Let  $\varepsilon > 0$  be given. Then by the triangle inequality,

$$\begin{aligned} |fg(y) - LK| &\leq |fg(y) - f(y)K| + |f(y)K - LK| \\ &\leq |f(y)| |g(y) - K| + |K| |f(y) - L|. \end{aligned} \quad (5.12)$$

There exists  $\delta_1$  such that if  $0 < |y - x| < \delta_1$ , then

$$|f(y) - L| < 1,$$

and so for such  $y$ , and the triangle inequality,  $|f(y)| < 1 + |L|$ . Therefore, for  $0 < |y - x| < \delta_1$ ,

$$|fg(y) - LK| \leq (1 + |K| + |L|) [|g(y) - K| + |f(y) - L|]. \quad (5.13)$$

Now let  $0 < \delta_2$  be such that for  $0 < |x - y| < \delta_2$ ,

$$|f(y) - L| < \frac{\varepsilon}{2(1 + |K| + |L|)}, \quad |g(y) - K| < \frac{\varepsilon}{2(1 + |K| + |L|)}.$$

Then letting  $0 < \delta \leq \min(\delta_1, \delta_2)$ , it follows from (5.13) that

$$|fg(y) - LK| < \varepsilon$$

and this proves (5.9). Limits as  $x \rightarrow \pm\infty$  and one sided limits are handled similarly.

The proof of (5.10) is left to you. It is just like the theorem about the quotient of continuous functions being continuous provided the function in the denominator is non zero at the point of interest.

Consider (5.11). Since  $h$  is continuous near  $L$ , it follows that for  $\varepsilon > 0$  given, there exists  $\eta > 0$  such that if  $|y-L| < \eta$ , then

$$|h(y) - h(L)| < \varepsilon$$

Now since  $\lim_{y \rightarrow x} f(y) = L$ , there exists  $\delta > 0$  such that if  $0 < |y-x| < \delta$ , then

$$|f(y) - L| < \eta.$$

Therefore, if  $0 < |y-x| < \delta$ ,

$$|h(f(y)) - h(L)| < \varepsilon.$$

The same theorem holds for one sided limits and limits as the variable moves toward  $\pm\infty$ . The proofs are left to you. They are minor modifications of the above.

It only remains to verify the last assertion. Assume  $f(y) \leq a$ . It is required to show that  $L \leq a$ . If this is not true, then  $L > a$ . Letting  $\varepsilon$  be small enough that  $a < L - \varepsilon$ , it follows that ultimately, for  $y$  close enough to  $x$ ,  $f(y) \in (L - \varepsilon, L + \varepsilon)$  which requires  $f(y) > a$  contrary to assumption.

A very useful theorem for finding limits is called the squeezing theorem.

**Theorem 5.9.5** Suppose  $\lim_{x \rightarrow a} f(x) = L = \lim_{x \rightarrow a} g(x)$  and for all  $x$  near  $a$ ,

$$f(x) \leq h(x) \leq g(x).$$

Then

$$\lim_{x \rightarrow a} h(x) = L.$$

**Proof:** If  $L \geq h(x)$ , then

$$|h(x) - L| \leq |f(x) - L|.$$

If  $L < h(x)$ , then

$$|h(x) - L| \leq |g(x) - L|.$$

Therefore,

$$|h(x) - L| \leq |f(x) - L| + |g(x) - L|.$$

Now let  $\varepsilon > 0$  be given. There exists  $\delta_1$  such that if  $0 < |x-a| < \delta_1$ ,

$$|f(x) - L| < \varepsilon/2$$

and there exists  $\delta_2$  such that if  $0 < |x-a| < \delta_2$ , then

$$|g(x) - L| < \varepsilon/2.$$

Letting  $0 < \delta \leq \min(\delta_1, \delta_2)$ , if  $0 < |x-a| < \delta$ , then

$$\begin{aligned} |h(x) - L| &\leq |f(x) - L| + |g(x) - L| \\ &< \varepsilon/2 + \varepsilon/2 = \varepsilon. \end{aligned}$$

This proves the theorem.

**Theorem 5.9.6** For  $f : I \rightarrow \mathbb{R}$ , and  $I$  is an interval of the form  $(a, b)$ ,  $[a, b)$ ,  $(a, b]$ , or  $[a, b]$ , then  $f$  is continuous at  $x \in I$  if and only if  $\lim_{y \rightarrow x} f(y) = f(x)$ .

**Proof:** You fill in the details. Compare the definition of continuous and the definition of the limit just given.

**Example 5.9.7** Find  $\lim_{x \rightarrow 3} \frac{x^2 - 9}{x - 3}$ .

Note that  $\frac{x^2 - 9}{x - 3} = x + 3$  whenever  $x \neq 3$ . Therefore, if  $0 < |x - 3| < \varepsilon$ ,

$$\left| \frac{x^2 - 9}{x - 3} - 6 \right| = |x + 3 - 6| = |x - 3| < \varepsilon.$$

It follows from the definition that this limit equals 6.

You should be careful to note that in the definition of limit, the variable **never equals the thing it is getting close to**. In this example,  $x$  is never equal to 3. This is very significant because, in interesting limits, the function whose limit is being taken will usually not be defined at the point of interest.

**Example 5.9.8** Let

$$f(x) = \frac{x^2 - 9}{x - 3} \text{ if } x \neq 3.$$

How should  $f$  be defined at  $x = 3$  so that the resulting function will be continuous there?

In the previous example, the limit of this function equals 6. Therefore, by Theorem 5.9.6 it is necessary to define  $f(3) \equiv 6$ .

**Example 5.9.9** Find  $\lim_{x \rightarrow \infty} \frac{x}{1+x}$ .

Write  $\frac{x}{1+x} = \frac{1}{1+(1/x)}$ . Now it seems clear that  $\lim_{x \rightarrow \infty} 1 + (1/x) = 1 \neq 0$ . Therefore, Theorem 5.9.4 implies

$$\lim_{x \rightarrow \infty} \frac{x}{1+x} = \lim_{x \rightarrow \infty} \frac{1}{1+(1/x)} = \frac{1}{1} = 1.$$

**Example 5.9.10** Show  $\lim_{x \rightarrow a} \sqrt{x} = \sqrt{a}$  whenever  $a \geq 0$ . In the case that  $a = 0$ , take the limit from the right.

There are two cases. First consider the case when  $a > 0$ . Let  $\varepsilon > 0$  be given. Multiply and divide by  $\sqrt{x} + \sqrt{a}$ . This yields

$$|\sqrt{x} - \sqrt{a}| = \left| \frac{x - a}{\sqrt{x} + \sqrt{a}} \right|.$$

Now let  $0 < \delta_1 < a/2$ . Then if  $|x - a| < \delta_1$ ,  $x > a/2$  and so

$$\begin{aligned} |\sqrt{x} - \sqrt{a}| &= \left| \frac{x - a}{\sqrt{x} + \sqrt{a}} \right| \leq \frac{|x - a|}{(\sqrt{a}/\sqrt{2}) + \sqrt{a}} \\ &\leq \frac{2\sqrt{2}}{\sqrt{a}} |x - a|. \end{aligned}$$

Now let  $0 < \delta \leq \min\left(\delta_1, \frac{\varepsilon\sqrt{a}}{2\sqrt{2}}\right)$ . Then for  $0 < |x - a| < \delta$ ,

$$|\sqrt{x} - \sqrt{a}| \leq \frac{2\sqrt{2}}{\sqrt{a}} |x - a| < \frac{2\sqrt{2}}{\sqrt{a}} \frac{\varepsilon\sqrt{a}}{2\sqrt{2}} = \varepsilon.$$

Next consider the case where  $a = 0$ . In this case, let  $\varepsilon > 0$  and let  $\delta = \varepsilon^2$ . Then if  $0 < x - 0 < \delta = \varepsilon^2$ , it follows that  $0 \leq \sqrt{x} < (\varepsilon^2)^{1/2} = \varepsilon$ .

## 5.10 Exercises

1. Find the following limits if possible

(a)  $\lim_{x \rightarrow 0+} \frac{|x|}{x}$

(b)  $\lim_{x \rightarrow 0+} \frac{x}{|x|}$

(c)  $\lim_{x \rightarrow 0-} \frac{|x|}{x}$

(d)  $\lim_{x \rightarrow 4} \frac{x^2-16}{x+4}$

(e)  $\lim_{x \rightarrow 3} \frac{x^2-9}{x+3}$

(f)  $\lim_{x \rightarrow -2} \frac{x^2-4}{x-2}$

(g)  $\lim_{x \rightarrow \infty} \frac{x}{1+x^2}$

(h)  $\lim_{x \rightarrow \infty} -2 \frac{x}{1+x^2}$

2. Find  $\lim_{h \rightarrow 0} \frac{\frac{1}{(x+h)^3} - \frac{1}{x^3}}{h}$ .

3. Find  $\lim_{x \rightarrow 4} \frac{\sqrt[4]{x} - \sqrt{2}}{\sqrt{x} - 2}$ .

4. Find  $\lim_{x \rightarrow \infty} \frac{\sqrt[5]{3x} + \sqrt[4]{x} + 7\sqrt{x}}{\sqrt{3x+1}}$ .

5. Find  $\lim_{x \rightarrow \infty} \frac{(x-3)^{20}(2x+1)^{30}}{(2x^2+7)^{25}}$ .

6. Find  $\lim_{x \rightarrow 2} \frac{x^2-4}{x^3+3x^2-9x-2}$ .

7. Find  $\lim_{x \rightarrow \infty} (\sqrt{1-7x+x^2} - \sqrt{1+7x+x^2})$ .

8. Prove Theorem 5.9.2 for right, left and limits as  $y \rightarrow \infty$ .

9. Prove from the definition that  $\lim_{x \rightarrow a} \sqrt[3]{x} = \sqrt[3]{a}$  for all  $a \in \mathbb{R}$ . **Hint:** You might want to use the formula for the difference of two cubes,

$$a^3 - b^3 = (a - b)(a^2 + ab + b^2).$$

10. Find  $\lim_{h \rightarrow 0} \frac{(x+h)^2 - x^2}{h}$ .

11. Prove Theorem 5.9.6 from the definitions of limit and continuity.

12. Find  $\lim_{h \rightarrow 0} \frac{(x+h)^3 - x^3}{h}$

13. Find  $\lim_{h \rightarrow 0} \frac{\frac{1}{x+h} - \frac{1}{x}}{h}$

14. Find  $\lim_{x \rightarrow -3} \frac{x^3+27}{x+3}$

15. Find  $\lim_{h \rightarrow 0} \frac{\sqrt{(3+h)^2} - 3}{h}$  if it exists.

16. Find the values of  $x$  for which  $\lim_{h \rightarrow 0} \frac{\sqrt{(x+h)^2} - x}{h}$  exists and find the limit.

17. Find  $\lim_{h \rightarrow 0} \frac{\sqrt[3]{(x+h)} - \sqrt[3]{x}}{h}$  if it exists. Here  $x \neq 0$ .

18. Suppose  $\lim_{y \rightarrow x+} f(y) = L_1 \neq L_2 = \lim_{y \rightarrow x-} f(y)$ . Show  $\lim_{y \rightarrow x} f(y)$  does not exist. **Hint:** Roughly, the argument goes as follows: For  $|y_1 - x|$  small and  $y_1 > x$ ,  $|f(y_1) - L_1|$  is small. Also, for  $|y_2 - x|$  small and  $y_2 < x$ ,  $|f(y_2) - L_2|$  is small. However, if a limit existed, then  $f(y_2)$  and  $f(y_1)$  would both need to be close to some number and so both  $L_1$  and  $L_2$  would need to be close to some number. However, this is impossible because they are different.
19. Show  $\lim_{x \rightarrow 0} \frac{\sin x}{x} = 1$ . **Hint:** You might consider Theorem 5.5.1 on Page 92 to write the inequality  $|\sin x| + 1 - \cos x \geq |x| \geq |\sin x|$  whenever  $|x|$  is small. Then divide both sides by  $|\sin x|$  and use some trig. identities to write  $\frac{\sin^2 x}{|\sin x|(1 + \cos x)} + 1 \geq \frac{|x|}{|\sin x|} \geq 1$  and then use squeezing theorem.

## 5.11 The Limit Of A Sequence

A closely related concept is the limit of a sequence. This was defined precisely a little before the definition of the limit by Bolzano<sup>4</sup>. The following is the precise definition of what is meant by the limit of a sequence.

**Definition 5.11.1** A sequence  $\{a_n\}_{n=1}^{\infty}$  converges to  $a$ ,

$$\lim_{n \rightarrow \infty} a_n = a \text{ or } a_n \rightarrow a$$

if and only if for every  $\varepsilon > 0$  there exists  $n_\varepsilon$  such that whenever  $n \geq n_\varepsilon$ ,

$$|a_n - a| < \varepsilon.$$

In words the definition says that given any measure of closeness,  $\varepsilon$ , the terms of the sequence are eventually all this close to  $a$ . Note the similarity with the concept of limit. Here, the word “eventually” refers to  $n$  being sufficiently large. Earlier, it referred to  $y$  being sufficiently close to  $x$  on one side or another or else  $x$  being sufficiently large in either the positive or negative directions.

**Theorem 5.11.2** If  $\lim_{n \rightarrow \infty} a_n = a$  and  $\lim_{n \rightarrow \infty} a_n = a_1$  then  $a_1 = a$ .

**Proof:** Suppose  $a_1 \neq a$ . Then let  $0 < \varepsilon < |a_1 - a|/2$  in the definition of the limit. It follows there exists  $n_\varepsilon$  such that if  $n \geq n_\varepsilon$ , then  $|a_n - a| < \varepsilon$  and  $|a_n - a_1| < \varepsilon$ . Therefore, for such  $n$ ,

$$\begin{aligned} |a_1 - a| &\leq |a_1 - a_n| + |a_n - a| \\ &< \varepsilon + \varepsilon < |a_1 - a|/2 + |a_1 - a|/2 = |a_1 - a|, \end{aligned}$$

a contradiction.

**Example 5.11.3** Let  $a_n = \frac{1}{n^2+1}$ .

Then it seems clear that

$$\lim_{n \rightarrow \infty} \frac{1}{n^2+1} = 0.$$

---

<sup>4</sup>Bernhard Bolzano lived from 1781 to 1848. He was a Catholic priest and held a position in philosophy at the University of Prague. He had strong views about the absurdity of war, educational reform, and the need for individual conscience. His convictions got him in trouble with Emperor Franz I of Austria and when he refused to recant, was forced out of the university. He understood the need for absolute rigor in mathematics. He also did work on physics.

In fact, this is true from the definition. Let  $\varepsilon > 0$  be given. Let  $n_\varepsilon \geq \sqrt{\varepsilon^{-1}}$ . Then if

$$n > n_\varepsilon \geq \sqrt{\varepsilon^{-1}},$$

it follows that  $n^2 + 1 > \varepsilon^{-1}$  and so

$$0 < \frac{1}{n^2 + 1} = a_n < \varepsilon.$$

Thus  $|a_n - 0| < \varepsilon$  whenever  $n$  is this large.

Note the definition was of no use in finding a candidate for the limit. This had to be produced based on other considerations. The definition is for verifying beyond any doubt that something is the limit. It is also what must be referred to in establishing theorems which are good for finding limits.

**Example 5.11.4** Let  $a_n = n^2$

Then in this case  $\lim_{n \rightarrow \infty} a_n$  does not exist. Sometimes this situation is also referred to by saying  $\lim_{n \rightarrow \infty} a_n = \infty$ .

**Example 5.11.5** Let  $a_n = (-1)^n$ .

In this case,  $\lim_{n \rightarrow \infty} (-1)^n$  does not exist. This follows from the definition. Let  $\varepsilon = 1/2$ . If there exists a limit,  $l$ , then eventually, for all  $n$  large enough,  $|a_n - l| < 1/2$ . However,  $|a_n - a_{n+1}| = 2$  and so,

$$2 = |a_n - a_{n+1}| \leq |a_n - l| + |l - a_{n+1}| < 1/2 + 1/2 = 1$$

which cannot hold. Therefore, there can be no limit for this sequence.

**Theorem 5.11.6** Suppose  $\{a_n\}$  and  $\{b_n\}$  are sequences and that

$$\lim_{n \rightarrow \infty} a_n = a \text{ and } \lim_{n \rightarrow \infty} b_n = b.$$

Also suppose  $x$  and  $y$  are real numbers. Then

$$\lim_{n \rightarrow \infty} xa_n + yb_n = xa + yb \quad (5.14)$$

$$\lim_{n \rightarrow \infty} a_n b_n = ab \quad (5.15)$$

If  $b \neq 0$ ,

$$\lim_{n \rightarrow \infty} \frac{a_n}{b_n} = \frac{a}{b}. \quad (5.16)$$

**Proof:** The first of these claims is left for you to do. To do the second, let  $\varepsilon > 0$  be given and choose  $n_1$  such that if  $n \geq n_1$  then

$$|a_n - a| < 1.$$

Then for such  $n$ , the triangle inequality implies

$$\begin{aligned} |a_n b_n - ab| &\leq |a_n b_n - a_n b| + |a_n b - ab| \\ &\leq |a_n| |b_n - b| + |b| |a_n - a| \\ &\leq (|a| + 1) |b_n - b| + |b| |a_n - a|. \end{aligned}$$



Now let  $n_2$  be large enough that for  $n \geq n_2$ ,

$$|b_n - b| < \frac{\varepsilon}{2(|a| + 1)}, \text{ and } |a_n - a| < \frac{\varepsilon}{2(|b| + 1)}.$$

Such a number exists because of the definition of limit. Therefore, let

$$n_\varepsilon > \max(n_1, n_2).$$

For  $n \geq n_\varepsilon$ ,

$$\begin{aligned} |a_n b_n - ab| &\leq (|a| + 1)|b_n - b| + |b||a_n - a| \\ &< (|a| + 1) \frac{\varepsilon}{2(|a| + 1)} + |b| \frac{\varepsilon}{2(|b| + 1)} \leq \varepsilon. \end{aligned}$$

This proves (5.15). Next consider (5.16).

Let  $\varepsilon > 0$  be given and let  $n_1$  be so large that whenever  $n \geq n_1$ ,

$$|b_n - b| < \frac{|b|}{2}.$$

Thus for such  $n$ ,

$$\begin{aligned} \left| \frac{a_n}{b_n} - \frac{a}{b} \right| &= \left| \frac{a_n b - ab_n}{bb_n} \right| \leq \frac{2}{|b|^2} [|a_n b - ab| + |ab - ab_n|] \\ &\leq \frac{2}{|b|} |a_n - a| + \frac{2|a|}{|b|^2} |b_n - b|. \end{aligned}$$

Now choose  $n_2$  so large that if  $n \geq n_2$ , then

$$|a_n - a| < \frac{\varepsilon |b|}{4}, \text{ and } |b_n - b| < \frac{\varepsilon |b|^2}{4(|a| + 1)}.$$

Letting  $n_\varepsilon > \max(n_1, n_2)$ , it follows that for  $n \geq n_\varepsilon$ ,

$$\begin{aligned} \left| \frac{a_n}{b_n} - \frac{a}{b} \right| &\leq \frac{2}{|b|} |a_n - a| + \frac{2|a|}{|b|^2} |b_n - b| \\ &< \frac{2}{|b|} \frac{\varepsilon |b|}{4} + \frac{2|a|}{|b|^2} \frac{\varepsilon |b|^2}{4(|a| + 1)} < \varepsilon. \end{aligned}$$

Another very useful theorem for finding limits is the squeezing theorem.

**Theorem 5.11.7** Suppose  $\lim_{n \rightarrow \infty} a_n = a = \lim_{n \rightarrow \infty} b_n$  and  $a_n \leq c_n \leq b_n$  for all  $n$  large enough. Then  $\lim_{n \rightarrow \infty} c_n = a$ .

**Proof:** Let  $\varepsilon > 0$  be given and let  $n_1$  be large enough that if  $n \geq n_1$ ,

$$|a_n - a| < \varepsilon/2 \text{ and } |b_n - a| < \varepsilon/2.$$

Then for such  $n$ ,

$$|c_n - a| \leq |a_n - a| + |b_n - a| < \varepsilon.$$

This proves the theorem.

As an example, consider the following.

**Example 5.11.8** *Let*

$$c_n \equiv (-1)^n \frac{1}{n}$$

*and let  $b_n = \frac{1}{n}$ , and  $a_n = -\frac{1}{n}$ . Then you may easily show that*

$$\lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} b_n = 0.$$

*Since  $a_n \leq c_n \leq b_n$ , it follows  $\lim_{n \rightarrow \infty} c_n = 0$  also.*

**Theorem 5.11.9**  $\lim_{n \rightarrow \infty} r^n = 0$ . *Whenever  $|r| < 1$ .*

**Proof:** If  $0 < r < 1$  it follows  $r^{-1} > 1$ . Why? Letting  $\alpha = \frac{1}{r} - 1$ , it follows

$$r = \frac{1}{1 + \alpha}.$$

Therefore, by the binomial theorem,

$$0 < r^n = \frac{1}{(1 + \alpha)^n} \leq \frac{1}{1 + \alpha n}.$$

Therefore,  $\lim_{n \rightarrow \infty} r^n = 0$  if  $0 < r < 1$ . Now in general, if  $|r| < 1$ ,  $|r^n| = |r|^n \rightarrow 0$  by the first part. This proves the theorem.

An important theorem is the one which states that if a sequence converges, so does every subsequence. You should review Definition 5.1.18 on Page 85 at this point.

**Theorem 5.11.10** *Let  $\{x_n\}$  be a sequence with  $\lim_{n \rightarrow \infty} x_n = x$  and let  $\{x_{n_k}\}$  be a subsequence. Then  $\lim_{k \rightarrow \infty} x_{n_k} = x$ .*

**Proof:** Let  $\varepsilon > 0$  be given. Then there exists  $n_\varepsilon$  such that if  $n > n_\varepsilon$ , then  $|x_n - x| < \varepsilon$ . Suppose  $k > n_\varepsilon$ . Then  $n_k \geq k > n_\varepsilon$  and so

$$|x_{n_k} - x| < \varepsilon$$

showing  $\lim_{k \rightarrow \infty} x_{n_k} = x$  as claimed.

### 5.11.1 Sequences And Completeness

You recall the definition of completeness which stated that every nonempty set of real numbers which is bounded above has a least upper bound and that every nonempty set of real numbers which is bounded below has a greatest lower bound and this is a property of the real line known as the completeness axiom. Geometrically, this involved filling in the holes. There is another way of describing completeness in terms of sequences which I believe is more useful than the least upper bound and greatest lower bound property.

**Definition 5.11.11**  $\{a_n\}$  *is a Cauchy sequence if for all  $\varepsilon > 0$ , there exists  $n_\varepsilon$  such that whenever  $n, m \geq n_\varepsilon$ ,*

$$|a_n - a_m| < \varepsilon.$$

A sequence is Cauchy means the terms are “bunching up to each other” as  $m, n$  get large.

**Theorem 5.11.12** *The set of terms in a Cauchy sequence in  $\mathbb{R}$  is bounded above and below.*

**Proof:** Let  $\varepsilon = 1$  in the definition of a Cauchy sequence and let  $n > n_1$ . Then from the definition,

$$|a_n - a_{n_1}| < 1.$$

It follows that for all  $n > n_1$ ,

$$|a_n| < 1 + |a_{n_1}|.$$

Therefore, for all  $n$ ,

$$|a_n| \leq 1 + |a_{n_1}| + \sum_{k=1}^{n_1} |a_k|.$$

This proves the theorem.

**Theorem 5.11.13** *If a sequence  $\{a_n\}$  in  $\mathbb{R}$  converges, then the sequence is a Cauchy sequence.*

**Proof:** Let  $\varepsilon > 0$  be given and suppose  $a_n \rightarrow a$ . Then from the definition of convergence, there exists  $n_\varepsilon$  such that if  $n > n_\varepsilon$ , it follows that

$$|a_n - a| < \frac{\varepsilon}{2}$$

Therefore, if  $m, n \geq n_\varepsilon + 1$ , it follows that

$$|a_n - a_m| \leq |a_n - a| + |a - a_m| < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon$$

showing that, since  $\varepsilon > 0$  is arbitrary,  $\{a_n\}$  is a Cauchy sequence.

**Definition 5.11.14** *The sequence,  $\{a_n\}$ , is monotone increasing if for all  $n$ ,  $a_n \leq a_{n+1}$ . The sequence is monotone decreasing if for all  $n$ ,  $a_n \geq a_{n+1}$ .*

If someone says a sequence is monotone, it usually means monotone increasing. There exists different descriptions of the completeness axiom. If you like you can simply add the three new criteria in the following theorem to the list of things which you mean when you say  $\mathbb{R}$  is complete and skip the proof. All versions of completeness involve the notion of filling in holes and they are really just different ways of expressing this idea.

In practice, it is often more convenient to use the first of the three equivalent versions of completeness in the following theorem which states that every Cauchy sequence converges. In fact, this version of completeness, although it is equivalent to the completeness axiom for the real line, also makes sense in many situations where Definition 2.14.1 on Page 42 does not make sense. For example, the concept of completeness is often needed in settings where there is no order. This happens as soon as one does multivariable calculus. From now on completeness will mean any of the three conditions in the following theorem.

It is the concept of completeness and the notion of limits which sets analysis apart from algebra. You will find that every existence theorem, a theorem which asserts the existence of something, in analysis depends on the assumption that some space is complete.

**Theorem 5.11.15** *The following conditions are equivalent to completeness.*

1. *Every Cauchy sequence converges*
2. *Every monotone increasing sequence which is bounded above converges.*
3. *Every monotone decreasing sequence which is bounded below converges.*

**Proof:** Suppose every Cauchy sequence converges and let  $S$  be a non empty set which is bounded above. In what follows,  $s_n \in S$  and  $b_n$  will be an upper bound of  $S$ . If, in the process about to be described,  $s_n = b_n$ , this will have shown the existence of a least upper bound to  $S$ . Therefore, assume  $s_n < b_n$  for all  $n$ . Let  $b_1$  be an upper bound of  $S$  and let  $s_1$  be an element of  $S$ . Suppose  $s_1, \dots, s_n$  and  $b_1, \dots, b_n$  have been chosen such that  $s_k \leq s_{k+1}$  and  $b_k \geq b_{k+1}$ . Consider  $\frac{s_n + b_n}{2}$ , the point on  $\mathbb{R}$  which is mid way between  $s_n$  and  $b_n$ . If this point is an upper bound, let

$$b_{n+1} \equiv \frac{s_n + b_n}{2}$$

and  $s_{n+1} = s_n$ . If the point is not an upper bound, let

$$s_{n+1} \in \left( \frac{s_n + b_n}{2}, b_n \right)$$

and let  $b_{n+1} = b_n$ . It follows this specifies an increasing sequence  $\{s_n\}$  and a decreasing sequence  $\{b_n\}$  such that

$$0 \leq b_n - s_n \leq 2^{-n} (b_1 - s_1).$$

Now if  $n > m$ ,

$$\begin{aligned} 0 \leq b_m - b_n &= |b_m - b_n| \\ &= \sum_{k=m}^{n-1} b_k - b_{k+1} \leq \sum_{k=m}^{n-1} b_k - s_k \leq \sum_{k=m}^{n-1} 2^{-k} (b_1 - s_1) \\ &= \frac{2^{-m} - 2^{-n}}{2^{-1}} (b_1 - s_1) \leq 2^{-m+1} (b_1 - s_1) \end{aligned}$$

and  $\lim_{m \rightarrow \infty} 2^{-m} = 0$  by Theorem 5.11.9. Therefore,  $\{b_n\}$  is a Cauchy sequence. Similarly,  $\{s_n\}$  is a Cauchy sequence. Let  $l \equiv \lim_{n \rightarrow \infty} s_n$  and let  $l_1 \equiv \lim_{n \rightarrow \infty} b_n$ . If  $n$  is large enough,

$$|l - s_n| < \varepsilon/3, |l_1 - b_n| < \varepsilon/3, \text{ and } |b_n - s_n| < \varepsilon/3.$$

Then

$$\begin{aligned} |l - l_1| &\leq |l - s_n| + |s_n - b_n| + |b_n - l_1| \\ &< \varepsilon/3 + \varepsilon/3 + \varepsilon/3 = \varepsilon. \end{aligned}$$

Since  $\varepsilon > 0$  is arbitrary,  $l = l_1$ . Why? Then  $l$  must be the least upper bound of  $S$ . It is an upper bound because if there were  $s > l$  where  $s \in S$ , then by the definition of limit,  $b_n < s$  for some  $n$ , violating the assumption that each  $b_n$  is an upper bound for  $S$ . On the other hand, if  $l_0 < l$ , then for all  $n$  large enough,  $s_n > l_0$ , which implies  $l_0$  is not an upper bound. This shows 1 implies completeness.

First note that 2 and 3 are equivalent. Why? Suppose 2 and consequently 3. Then the same construction yields the two monotone sequences, one increasing and the other decreasing. The sequence  $\{b_n\}$  is bounded below by  $s_m$  for all  $m$  and the sequence  $\{s_n\}$  is bounded above by  $b_m$  for all  $m$ . Why? Therefore, the two sequences converge. The rest of the argument is the same as the above. Thus 2 and 3 imply completeness.

Now suppose completeness and let  $\{a_n\}$  be an increasing sequence which is bounded above. Let  $a$  be the least upper bound of the set of points in the sequence. If  $\varepsilon > 0$  is given, there exists  $n_\varepsilon$  such that  $a - \varepsilon < a_{n_\varepsilon}$ . Since  $\{a_n\}$  is a monotone sequence, it follows that whenever  $n > n_\varepsilon$ ,  $a - \varepsilon < a_n \leq a$ . This proves  $\lim_{n \rightarrow \infty} a_n = a$  and proves convergence. Since 3 is equivalent to 2, this is also established. It follows 3 and 2 are equivalent to completeness. It remains to show that completeness implies every Cauchy sequence converges.

Suppose completeness and let  $\{a_n\}$  be a Cauchy sequence. Let

$$\inf \{a_k : k \geq n\} \equiv A_n, \sup \{a_k : k \geq n\} \equiv B_n$$

Then  $A_n$  is an increasing sequence while  $B_n$  is a decreasing sequence and  $B_n \geq A_n$ . Furthermore,

$$\lim_{n \rightarrow \infty} B_n - A_n = 0.$$

The details of these assertions are easy and are left to the reader. Also,  $\{A_n\}$  is bounded below by any lower bound for the original Cauchy sequence while  $\{B_n\}$  is bounded above by any upper bound for the original Cauchy sequence. By the equivalence of completeness with 3 and 2, it follows there exists  $a$  such that  $a = \lim_{n \rightarrow \infty} A_n = \lim_{n \rightarrow \infty} B_n$ . Since  $B_n \geq a_n \geq A_n$ , the squeezing theorem implies  $\lim_{n \rightarrow \infty} a_n = a$  and this proves the equivalence of these characterizations of completeness.

**Theorem 5.11.16** *Let  $\{a_n\}$  be a monotone increasing sequence which is bounded above. Then  $\lim_{n \rightarrow \infty} a_n = \sup \{a_n : n \geq 1\}$*

**Proof:** Let  $a = \sup \{a_n : n \geq 1\}$  and let  $\varepsilon > 0$  be given. Then from Proposition 2.14.3 on Page 42 there exists  $m$  such that  $a - \varepsilon < a_m \leq a$ . Since the sequence is increasing, it follows that for all  $n \geq m$ ,  $a - \varepsilon < a_n \leq a$ . Thus  $a = \lim_{n \rightarrow \infty} a_n$ .

### 5.11.2 Decimals

You are all familiar with decimals. In the United States these are written in the form  $.a_1a_2a_3\cdots$  where the  $a_i$  are integers between 0 and 9.<sup>5</sup> Thus .23417432 is a number written as a decimal. You also recall the meaning of such notation in the case of a terminating decimal. For example, .234 is defined as  $\frac{2}{10} + \frac{3}{10^2} + \frac{4}{10^3}$ . Now what is meant by a nonterminating decimal?

**Definition 5.11.17** *Let  $.a_1a_2\cdots$  be a decimal. Define*

$$.a_1a_2\cdots \equiv \lim_{n \rightarrow \infty} \sum_{k=1}^n \frac{a_k}{10^k}.$$

**Proposition 5.11.18** *The above definition makes sense.*

**Proof:** Note the sequence  $\{\sum_{k=1}^n \frac{a_k}{10^k}\}_{n=1}^\infty$  is an increasing sequence. Therefore, if there exists an upper bound, it follows from Theorem 5.11.16 that this sequence converges and so the definition is well defined.

$$\sum_{k=1}^n \frac{a_k}{10^k} \leq \sum_{k=1}^n \frac{9}{10^k} = 9 \sum_{k=1}^n \frac{1}{10^k}.$$

Now

$$\begin{aligned} \frac{9}{10} \left( \sum_{k=1}^n \frac{1}{10^k} \right) &= \sum_{k=1}^n \frac{1}{10^k} - \frac{1}{10} \sum_{k=1}^n \frac{1}{10^k} = \sum_{k=1}^n \frac{1}{10^k} - \sum_{k=2}^{n+1} \frac{1}{10^k} \\ &= \frac{1}{10} - \frac{1}{10^{n+1}} \end{aligned}$$

---

<sup>5</sup>In France and Russia they use a comma instead of a period. This looks very strange but that is just the way they do it.

and so

$$\sum_{k=1}^n \frac{1}{10^k} \leq \frac{10}{9} \left( \frac{1}{10} - \frac{1}{10^{n+1}} \right) \leq \frac{10}{9} \left( \frac{1}{10} \right) = \frac{1}{9}.$$

Therefore, since this holds for all  $n$ , it follows the above sequence is bounded above. It follows the limit exists.

### 5.11.3 Continuity And The Limit Of A Sequence

There is a very useful way of thinking of continuity in terms of limits of sequences found in the following theorem. In words, it says a function is continuous if it takes convergent sequences to convergent sequences whenever possible.

**Theorem 5.11.19** *A function  $f : D(f) \rightarrow \mathbb{R}$  is continuous at  $x \in D(f)$  if and only if, whenever  $x_n \rightarrow x$  with  $x_n \in D(f)$ , it follows  $f(x_n) \rightarrow f(x)$ .*

**Proof:** Suppose first that  $f$  is continuous at  $x$  and let  $x_n \rightarrow x$ . Let  $\varepsilon > 0$  be given. By continuity, there exists  $\delta > 0$  such that if  $|y - x| < \delta$ , then  $|f(y) - f(x)| < \varepsilon$ . However, there exists  $n_\delta$  such that if  $n \geq n_\delta$ , then  $|x_n - x| < \delta$  and so for all  $n$  this large,

$$|f(x) - f(x_n)| < \varepsilon$$

which shows  $f(x_n) \rightarrow f(x)$ .

Now suppose the condition about taking convergent sequences to convergent sequences holds at  $x$ . Suppose  $f$  fails to be continuous at  $x$ . Then there exists  $\varepsilon > 0$  and  $x_n \in D(f)$  such that  $|x - x_n| < \frac{1}{n}$ , yet

$$|f(x) - f(x_n)| \geq \varepsilon.$$

But this is clearly a contradiction because, although  $x_n \rightarrow x$ ,  $f(x_n)$  fails to converge to  $f(x)$ . It follows  $f$  must be continuous after all. This proves the theorem.

## 5.12 Exercises

1. Find  $\lim_{n \rightarrow \infty} \frac{n}{3n+4}$ .
2. Find  $\lim_{n \rightarrow \infty} \frac{3n^4+7n+1000}{n^4+1}$ .
3. Find  $\lim_{n \rightarrow \infty} \frac{2^n+7(5^n)}{4^n+2(5^n)}$ .
4. Find  $\lim_{n \rightarrow \infty} n \tan \frac{1}{n}$ . **Hint:** See Problem 19 on Page 103.
5. Find  $\lim_{n \rightarrow \infty} n \sin \frac{2}{n}$ . **Hint:** See Problem 19 on Page 103.
6. Find  $\lim_{n \rightarrow \infty} \sqrt{(n \sin \frac{9}{n})}$ . **Hint:** See Problem 19 on Page 103.
7. Find  $\lim_{n \rightarrow \infty} \sqrt{(n^2 + 6n)} - n$ . **Hint:** Multiply and divide by  $\sqrt{(n^2 + 6n)} + n$ .
8. Find  $\lim_{n \rightarrow \infty} \sum_{k=1}^n \frac{1}{10^k}$ .
9. For  $|r| < 1$ , find  $\lim_{n \rightarrow \infty} \sum_{k=0}^n r^k$ . **Hint:** First show  $\sum_{k=0}^n r^k = \frac{r^{n+1}}{r-1} - \frac{1}{r-1}$ . Then recall Theorem 5.11.9.

10. Suppose  $x = .343434\overline{34}$  where the bar over the last 34 signifies that this repeats forever. In elementary school you were probably given the following procedure for finding the number  $x$  as a quotient of integers. First multiply by 100 to get  $100x = 34.3434\overline{34}$  and then subtract to get  $99x = 34$ . From this you conclude that  $x = 34/99$ . Fully justify this procedure. **Hint:**  $.343434\overline{34} = \lim_{n \rightarrow \infty} 34 \sum_{k=1}^n \left(\frac{1}{100}\right)^k$  now use Problem 9.
11. Suppose  $D(f) = [0, 1] \cup \{9\}$  and  $f(x) = x$  on  $[0, 1]$  while  $f(9) = 5$ . Is  $f$  continuous at the point, 9? Use whichever definition of continuity you like.
12. Suppose  $x_n \rightarrow x$  and  $x_n \leq c$ . Show that  $x \leq c$ . Also show that if  $x_n \rightarrow x$  and  $x_n \geq c$ , then  $x \geq c$ . **Hint:** If this is not true, argue that for all  $n$  large enough  $x_n > c$ .
13. Let  $a \in [0, 1]$ . Show  $a = .a_1a_2a_3 \cdots$  for a unique choice of integers,  $a_1, a_2, \dots$  if it is possible to do this. Otherwise, give an example.
14. Show every rational number between 0 and 1 has a decimal expansion which either repeats or terminates.
15. Consider the number whose decimal expansion is  $.010010001000010000010000001 \cdots$ . Show this is an irrational number. Now using this, show that between any two integers there exists an irrational number. Next show that between any two numbers there exists an irrational number.
16. Using the binomial theorem prove that for all  $n \in \mathbb{N}$ ,  $\left(1 + \frac{1}{n}\right)^n \leq \left(1 + \frac{1}{n+1}\right)^{n+1}$ .  
**Hint:** Show first that  $\binom{n}{k} = \frac{n \cdot (n-1) \cdots (n-k+1)}{k!}$ . By the binomial theorem,

$$\left(1 + \frac{1}{n}\right)^n = \sum_{k=0}^n \binom{n}{k} \left(\frac{1}{n}\right)^k = \sum_{k=0}^n \overbrace{\frac{n \cdot (n-1) \cdots (n-k+1)}{k! n^k}}^{k \text{ factors}}.$$

Now consider the term  $\frac{n \cdot (n-1) \cdots (n-k+1)}{k! n^k}$  and note that a similar term occurs in the binomial expansion for  $\left(1 + \frac{1}{n+1}\right)^{n+1}$  except you replace  $n$  with  $n+1$  wherever this occurs. Argue the term got bigger and then note that in the binomial expansion for  $\left(1 + \frac{1}{n+1}\right)^{n+1}$ , there are more terms.

17. Prove by induction that for all  $k \geq 4$ ,  $2^k \leq k!$
18. Use the Problems 21 and 16 to verify for all  $n \in \mathbb{N}$ ,  $\left(1 + \frac{1}{n}\right)^n \leq 3$ .
19. Prove  $\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n$  exists and equals a number less than 3.
20. Using Problem 18, prove  $n^{n+1} \geq (n+1)^n$  for all integers,  $n \geq 3$ .
21. Find  $\lim_{n \rightarrow \infty} n \sin n$  if it exists. If it does not exist, explain why it does not.
22. Recall the axiom of completeness states that a set which is bounded above has a least upper bound and a set which is bounded below has a greatest lower bound. Show that a monotone decreasing sequence which is bounded below converges to its greatest lower bound. **Hint:** Let  $a$  denote the greatest lower bound and recall that because of this, it follows that for all  $\varepsilon > 0$  there exist points of  $\{a_n\}$  in  $[a, a + \varepsilon]$ .

23. Let  $A_n = \sum_{k=2}^n \frac{1}{k(k-1)}$  for  $n \geq 2$ . Show  $\lim_{n \rightarrow \infty} A_n$  exists. **Hint:** Show there exists an upper bound to the  $A_n$  as follows.

$$\begin{aligned} \sum_{k=2}^n \frac{1}{k(k-1)} &= \sum_{k=2}^n \left( \frac{1}{k-1} - \frac{1}{k} \right) \\ &= \frac{1}{2} - \frac{1}{n-1} \leq \frac{1}{2}. \end{aligned}$$

24. Let  $H_n = \sum_{k=1}^n \frac{1}{k^2}$  for  $n \geq 2$ . Show  $\lim_{n \rightarrow \infty} H_n$  exists. **Hint:** Use the above problem to obtain the existence of an upper bound.
25. Let  $a$  be a positive number and let  $x_1 = b > 0$  where  $b^2 > a$ . Explain why there exists such a number,  $b$ . Now having defined  $x_n$ , define  $x_{n+1} \equiv \frac{1}{2} \left( x_n + \frac{a}{x_n} \right)$ . Verify that  $\{x_n\}$  is a decreasing sequence and that it satisfies  $x_n^2 \geq a$  for all  $n$  and is therefore, bounded below. Explain why  $\lim_{n \rightarrow \infty} x_n$  exists. If  $x$  is this limit, show that  $x^2 = a$ . Explain how this shows that every positive real number has a square root. This is an example of a recursively defined sequence. Note this does not give a formula for  $x_n$ , just a rule which tells us how to define  $x_{n+1}$  if  $x_n$  is known.
26. Let  $a_1 = 0$  and suppose that  $a_{n+1} = \frac{9}{9-a_n}$ . Write  $a_2, a_3, a_4$ . Now prove that for all  $n$ , it follows that  $a_n \leq \frac{9}{2} + \frac{3}{2}\sqrt{5}$  (By Problem 6 on Page 97 there is no problem with the existence of various roots of positive numbers.) and so the sequence is bounded above. Next show that the sequence is increasing and so it converges. Find the limit of the sequence. **Hint:** You should prove these things by induction. Finally, to find the limit, let  $n \rightarrow \infty$  in both sides and argue that the limit,  $a$ , must satisfy  $a = \frac{9}{9-a}$ .
27. If  $x \in \mathbb{R}$ , show there exists a sequence of rational numbers,  $\{x_n\}$  such that  $x_n \rightarrow x$  and a sequence of irrational numbers,  $\{x'_n\}$  such that  $x'_n \rightarrow x$ . Now consider the following function.

$$f(x) = \begin{cases} 1 & \text{if } x \text{ is rational} \\ 0 & \text{if } x \text{ is irrational} \end{cases}.$$

Show using the sequential version of continuity in Theorem 5.11.19 that  $f$  is discontinuous at every point.

28. If  $x \in \mathbb{R}$ , show there exists a sequence of rational numbers,  $\{x_n\}$  such that  $x_n \rightarrow x$  and a sequence of irrational numbers,  $\{x'_n\}$  such that  $x'_n \rightarrow x$ . Now consider the following function.

$$f(x) = \begin{cases} x & \text{if } x \text{ is rational} \\ 0 & \text{if } x \text{ is irrational} \end{cases}.$$

Show using the sequential version of continuity in Theorem 5.11.19 that  $f$  is continuous at 0 and nowhere else.

29. The nested interval lemma and Theorem 5.11.19 can be used to give an easy proof of the intermediate value theorem. Suppose  $f(a) > 0$  and  $f(b) < 0$  for  $f$  a continuous function defined on  $[a, b]$ . The intermediate value theorem states that under these conditions, there exists  $x \in (a, b)$  such that  $f(x) = 0$ . Prove this theorem as follows: Let  $c = \frac{a+b}{2}$  and consider the intervals  $[a, c]$  and  $[c, b]$ . Show that on one of these intervals,  $f$  is nonnegative at one end and nonpositive at the other. Now consider that interval, divide it in half as was done for the original interval and argue that on one of these smaller intervals, the function has different signs at the two endpoints. Continue in this way. Next apply the nested interval lemma to get  $x$  in all these intervals and



argue there exist sequences,  $x_n \rightarrow x$  and  $y_n \rightarrow x$  such that  $f(x_n) < 0$  and  $f(y_n) > 0$ . By continuity, you can assume  $f(x_n) \rightarrow f(x)$  and  $f(y_n) \rightarrow f(x)$ . Show this requires that  $f(x) = 0$ .

30. If  $\lim_{n \rightarrow \infty} a_n = a$ , does it follow that  $\lim_{n \rightarrow \infty} |a_n| = |a|$ ? Prove or else give a counter example.

31. Show the following converge to 0.

(a)  $\frac{n^5}{1.01^n}$

(b)  $\frac{10^n}{n!}$

32. Prove  $\lim_{n \rightarrow \infty} \sqrt[n]{n} = 1$ . **Hint:** Let  $e_n \equiv \sqrt[n]{n} - 1$  so that  $(1 + e_n)^n = n$ . Now observe that  $e_n > 0$  and use the binomial theorem to conclude  $1 + ne_n + \frac{n(n-1)}{2}e_n^2 \leq n$ . This nice approach to establishing this limit using only elementary algebra is in Rudin [14].

## 5.13 Uniform Continuity

There is a theorem about the integral of a continuous function which requires the notion of uniform continuity. This is discussed in this section. Consider the function  $f(x) = \frac{1}{x}$  for  $x \in (0, 1)$ . This is a continuous function because, by Theorem 5.4.1, it is continuous at every point of  $(0, 1)$ . However, for a given  $\varepsilon > 0$ , the  $\delta$  needed in the  $\varepsilon, \delta$  definition of continuity becomes very small as  $x$  gets close to 0. The notion of uniform continuity involves being able to choose a single  $\delta$  which works on the whole domain of  $f$ . Here is the definition.

**Definition 5.13.1** Let  $f : D \subseteq \mathbb{R} \rightarrow \mathbb{R}$  be a function. Then  $f$  is uniformly continuous if for every  $\varepsilon > 0$ , there exists a  $\delta$  **depending only on**  $\varepsilon$  such that if  $|x - y| < \delta$  then  $|f(x) - f(y)| < \varepsilon$ .

It is an amazing fact that under certain conditions continuity implies uniform continuity.

**Definition 5.13.2** A set,  $K \subseteq \mathbb{R}$  is sequentially compact if whenever  $\{a_n\} \subseteq K$  is a sequence, there exists a subsequence,  $\{a_{n_k}\}$  such that this subsequence converges to a point of  $K$ .

The following theorem is part of the Heine Borel theorem.

**Theorem 5.13.3** Every closed interval,  $[a, b]$  is sequentially compact.

**Proof:** Let  $\{x_n\} \subseteq [a, b] \equiv I_0$ . Consider the two intervals  $[a, \frac{a+b}{2}]$  and  $[\frac{a+b}{2}, b]$  each of which has length  $(b - a)/2$ . At least one of these intervals contains  $x_n$  for infinitely many values of  $n$ . Call this interval  $I_1$ . Now do for  $I_1$  what was done for  $I_0$ . Split it in half and let  $I_2$  be the interval which contains  $x_n$  for infinitely many values of  $n$ . Continue this way obtaining a sequence of nested intervals  $I_0 \supseteq I_1 \supseteq I_2 \supseteq I_3 \cdots$  where the length of  $I_n$  is  $(b - a)/2^n$ . Now pick  $n_1$  such that  $x_{n_1} \in I_1$ ,  $n_2$  such that  $n_2 > n_1$  and  $x_{n_2} \in I_2$ ,  $n_3$  such that  $n_3 > n_2$  and  $x_{n_3} \in I_3$ , etc. (This can be done because in each case the intervals contained  $x_n$  for infinitely many values of  $n$ .) By the nested interval lemma there exists a point,  $c$  contained in all these intervals. Furthermore,

$$|x_{n_k} - c| < (b - a)2^{-k}$$

and so  $\lim_{k \rightarrow \infty} x_{n_k} = c \in [a, b]$ . This proves the theorem.

**Theorem 5.13.4** *Let  $f : K \rightarrow \mathbb{R}$  be continuous where  $K$  is a sequentially compact set in  $\mathbb{R}$ . Then  $f$  is uniformly continuous on  $K$ .*

**Proof:** If this is not true, there exists  $\varepsilon > 0$  such that for every  $\delta > 0$  there exists a pair of points,  $x_\delta$  and  $y_\delta$  such that even though  $|x_\delta - y_\delta| < \delta$ ,  $|f(x_\delta) - f(y_\delta)| \geq \varepsilon$ . Taking a succession of values for  $\delta$  equal to  $1, 1/2, 1/3, \dots$ , and letting the exceptional pair of points for  $\delta = 1/n$  be denoted by  $x_n$  and  $y_n$ ,

$$|x_n - y_n| < \frac{1}{n}, |f(x_n) - f(y_n)| \geq \varepsilon.$$

Now since  $K$  is sequentially compact, there exists a subsequence,  $\{x_{n_k}\}$  such that  $x_{n_k} \rightarrow z \in K$ . Now  $n_k \geq k$  and so

$$|x_{n_k} - y_{n_k}| < \frac{1}{k}.$$

Consequently,  $y_{n_k} \rightarrow z$  also. ( $x_{n_k}$  is like a person walking toward a certain point and  $y_{n_k}$  is like a dog on a leash which is constantly getting shorter. Obviously  $y_{n_k}$  must also move toward the point also. You should give a precise proof of what is needed here.) By continuity of  $f$  and Problem 12 on Page 111,

$$0 = |f(z) - f(z)| = \lim_{k \rightarrow \infty} |f(x_{n_k}) - f(y_{n_k})| \geq \varepsilon,$$

an obvious contradiction. Therefore, the theorem must be true.

The following corollary follows from this theorem and Theorem 5.13.3.

**Corollary 5.13.5** *Suppose  $I$  is a closed interval,  $I = [a, b]$  and  $f : I \rightarrow \mathbb{R}$  is continuous. Then  $f$  is uniformly continuous.*

## 5.14 Exercises

1. A function,  $f : D \subseteq \mathbb{R} \rightarrow \mathbb{R}$  is Lipschitz continuous or just Lipschitz for short if there exists a constant,  $K$  such that

$$|f(x) - f(y)| \leq K|x - y|$$

for all  $x, y \in D$ . Show every Lipschitz function is uniformly continuous.

2. If  $|x_n - y_n| \rightarrow 0$  and  $x_n \rightarrow z$ , show that  $y_n \rightarrow z$  also.
3. Consider  $f : (1, \infty) \rightarrow \mathbb{R}$  given by  $f(x) = \frac{1}{x}$ . Show  $f$  is uniformly continuous even though the set on which  $f$  is defined is not sequentially compact.
4. If  $f$  is uniformly continuous, does it follow that  $|f|$  is also uniformly continuous? If  $|f|$  is uniformly continuous does it follow that  $f$  is uniformly continuous? Answer the same questions with “uniformly continuous” replaced with “continuous”. Explain why.

## 5.15 Theorems About Continuous Functions

In this section, proofs of some theorems which have not been proved yet are given.

**Theorem 5.15.1** *The following assertions are valid*

1. The function,  $af + bg$  is continuous at  $x$  when  $f, g$  are continuous at  $x \in D(f) \cap D(g)$  and  $a, b \in \mathbb{R}$ .
2. If  $f$  and  $g$  are each real valued functions continuous at  $x$ , then  $fg$  is continuous at  $x$ . If, in addition to this,  $g(x) \neq 0$ , then  $f/g$  is continuous at  $x$ .
3. If  $f$  is continuous at  $x$ ,  $f(x) \in D(g) \subseteq \mathbb{R}$ , and  $g$  is continuous at  $f(x)$ , then  $g \circ f$  is continuous at  $x$ .
4. The function  $f: \mathbb{R} \rightarrow \mathbb{R}$ , given by  $f(x) = |x|$  is continuous.

**Proof:** First consider 1.) Let  $\varepsilon > 0$  be given. By assumption, there exist  $\delta_1 > 0$  such that whenever  $|x - y| < \delta_1$ , it follows  $|f(x) - f(y)| < \frac{\varepsilon}{2(|a|+|b|+1)}$  and there exists  $\delta_2 > 0$  such that whenever  $|x - y| < \delta_2$ , it follows that  $|g(x) - g(y)| < \frac{\varepsilon}{2(|a|+|b|+1)}$ . Then let  $0 < \delta \leq \min(\delta_1, \delta_2)$ . If  $|x - y| < \delta$ , then everything happens at once. Therefore, using the triangle inequality

$$\begin{aligned} & |af(x) + bf(x) - (ag(y) + bg(y))| \\ & \leq |a||f(x) - f(y)| + |b||g(x) - g(y)| \\ & < |a| \left( \frac{\varepsilon}{2(|a|+|b|+1)} \right) + |b| \left( \frac{\varepsilon}{2(|a|+|b|+1)} \right) < \varepsilon. \end{aligned}$$

Now consider 2.) There exists  $\delta_1 > 0$  such that if  $|y - x| < \delta_1$ , then  $|f(x) - f(y)| < 1$ . Therefore, for such  $y$ ,

$$|f(y)| < 1 + |f(x)|.$$

It follows that for such  $y$ ,

$$\begin{aligned} |fg(x) - fg(y)| & \leq |f(x)g(x) - g(x)f(y)| + |g(x)f(y) - f(y)g(y)| \\ & \leq |g(x)||f(x) - f(y)| + |f(y)||g(x) - g(y)| \\ & \leq (1 + |g(x)| + |f(y)|)[|g(x) - g(y)| + |f(x) - f(y)|]. \end{aligned}$$

Now let  $\varepsilon > 0$  be given. There exists  $\delta_2$  such that if  $|x - y| < \delta_2$ , then

$$|g(x) - g(y)| < \frac{\varepsilon}{2(1 + |g(x)| + |f(y)|)},$$

and there exists  $\delta_3$  such that if  $|x - y| < \delta_3$ , then

$$|f(x) - f(y)| < \frac{\varepsilon}{2(1 + |g(x)| + |f(y)|)}$$

Now let  $0 < \delta \leq \min(\delta_1, \delta_2, \delta_3)$ . Then if  $|x - y| < \delta$ , all the above hold at once and so

$$\begin{aligned} & |fg(x) - fg(y)| \leq \\ & (1 + |g(x)| + |f(y)|)[|g(x) - g(y)| + |f(x) - f(y)|] \\ & < (1 + |g(x)| + |f(y)|) \left( \frac{\varepsilon}{2(1 + |g(x)| + |f(y)|)} + \frac{\varepsilon}{2(1 + |g(x)| + |f(y)|)} \right) = \varepsilon. \end{aligned}$$

This proves the first part of 2.) To obtain the second part, let  $\delta_1$  be as described above and let  $\delta_0 > 0$  be such that for  $|x - y| < \delta_0$ ,

$$|g(x) - g(y)| < |g(x)|/2$$

and so by the triangle inequality,

$$-|g(x)|/2 \leq |g(y)| - |g(x)| \leq |g(x)|/2$$

which implies  $|g(y)| \geq |g(x)|/2$ , and  $|g(y)| < 3|g(x)|/2$ .

Then if  $|x-y| < \min(\delta_0, \delta_1)$ ,

$$\begin{aligned} \left| \frac{f(x)}{g(x)} - \frac{f(y)}{g(y)} \right| &= \left| \frac{f(x)g(y) - f(y)g(x)}{g(x)g(y)} \right| \\ &\leq \frac{|f(x)g(y) - f(y)g(x)|}{\left(\frac{|g(x)|^2}{2}\right)} \\ &= \frac{2|f(x)g(y) - f(y)g(x)|}{|g(x)|^2} \\ &\leq \frac{2}{|g(x)|^2} [|f(x)g(y) - f(y)g(y) + f(y)g(y) - f(y)g(x)|] \\ &\leq \frac{2}{|g(x)|^2} [|g(y)||f(x) - f(y)| + |f(y)||g(y) - g(x)|] \\ &\leq \frac{2}{|g(x)|^2} \left[ \frac{3}{2}|g(x)||f(x) - f(y)| + (1 + |f(x)|)|g(y) - g(x)| \right] \\ &\leq \frac{2}{|g(x)|^2} (1 + 2|f(x)| + 2|g(x)|) [|f(x) - f(y)| + |g(y) - g(x)|] \\ &\equiv M [|f(x) - f(y)| + |g(y) - g(x)|] \end{aligned}$$

where  $M$  is defined by

$$M \equiv \frac{2}{|g(x)|^2} (1 + 2|f(x)| + 2|g(x)|)$$

Now let  $\delta_2$  be such that if  $|x-y| < \delta_2$ , then

$$|f(x) - f(y)| < \frac{\varepsilon}{2} M^{-1}$$

and let  $\delta_3$  be such that if  $|x-y| < \delta_3$ , then

$$|g(y) - g(x)| < \frac{\varepsilon}{2} M^{-1}.$$

Then if  $0 < \delta \leq \min(\delta_0, \delta_1, \delta_2, \delta_3)$ , and  $|x-y| < \delta$ , everything holds and

$$\begin{aligned} \left| \frac{f(x)}{g(x)} - \frac{f(y)}{g(y)} \right| &\leq M [|f(x) - f(y)| + |g(y) - g(x)|] \\ &< M \left[ \frac{\varepsilon}{2} M^{-1} + \frac{\varepsilon}{2} M^{-1} \right] = \varepsilon. \end{aligned}$$

This completes the proof of the second part of 2.)

Note that in these proofs no effort is made to find some sort of “best”  $\delta$ . The problem is one which has a yes or a no answer. Either is it or it is not continuous.

Now consider 3.). If  $f$  is continuous at  $x$ ,  $f(x) \in D(g) \subseteq \mathbb{R}^p$ , and  $g$  is continuous at  $f(x)$ , then  $g \circ f$  is continuous at  $x$ . Let  $\varepsilon > 0$  be given. Then there exists  $\eta > 0$  such that if

$|y - f(x)| < \eta$  and  $y \in D(g)$ , it follows that  $|g(y) - g(f(x))| < \varepsilon$ . From continuity of  $f$  at  $x$ , there exists  $\delta > 0$  such that if  $|x - z| < \delta$  and  $z \in D(f)$ , then  $|f(z) - f(x)| < \eta$ . Then if  $|x - z| < \delta$  and  $z \in D(g \circ f) \subseteq D(f)$ , all the above hold and so

$$|g(f(z)) - g(f(x))| < \varepsilon.$$

This proves part 3.)

To verify part 4.), let  $\varepsilon > 0$  be given and let  $\delta = \varepsilon$ . Then if  $|x - y| < \delta$ , the triangle inequality implies

$$\begin{aligned} |f(x) - f(y)| &= ||x| - |y|| \\ &\leq |x - y| < \delta = \varepsilon. \end{aligned}$$

This proves part 4.) and completes the proof of the theorem.

Next here is a proof of the intermediate value theorem.

**Theorem 5.15.2** *Suppose  $f : [a, b] \rightarrow \mathbb{R}$  is continuous and suppose  $f(a) < c < f(b)$ . Then there exists  $x \in (a, b)$  such that  $f(x) = c$ .*

**Proof:** Let  $d = \frac{a+b}{2}$  and consider the intervals  $[a, d]$  and  $[d, b]$ . If  $f(d) \geq c$ , then on  $[a, d]$ , the function is  $\leq c$  at one end point and  $\geq c$  at the other. On the other hand, if  $f(d) \leq c$ , then on  $[d, b]$   $f \geq 0$  at one end point and  $\leq 0$  at the other. Pick the interval on which  $f$  has values which are at least as large as  $c$  and values no larger than  $c$ . Now consider that interval, divide it in half as was done for the original interval and argue that on one of these smaller intervals, the function has values at least as large as  $c$  and values no larger than  $c$ . Continue in this way. Next apply the nested interval lemma to get  $x$  in all these intervals. In the  $n^{\text{th}}$  interval, let  $x_n, y_n$  be elements of this interval such that  $f(x_n) \leq c, f(y_n) \geq c$ . Now  $|x_n - x| \leq (b - a)2^{-n}$  and  $|y_n - x| \leq (b - a)2^{-n}$  and so  $x_n \rightarrow x$  and  $y_n \rightarrow x$ . Therefore,

$$f(x) - c = \lim_{n \rightarrow \infty} (f(x_n) - c) \leq 0$$

while

$$f(x) - c = \lim_{n \rightarrow \infty} (f(y_n) - c) \geq 0.$$

Consequently  $f(x) = c$  and this proves the theorem. (For the last step, see Problem 12 on Page 111).

**Lemma 5.15.3** *Let  $\phi : [a, b] \rightarrow \mathbb{R}$  be a continuous function and suppose  $\phi$  is 1-1 on  $(a, b)$ . Then  $\phi$  is either strictly increasing or strictly decreasing on  $[a, b]$ .*

**Proof:** First it is shown that  $\phi$  is either strictly increasing or strictly decreasing on  $(a, b)$ .

If  $\phi$  is not strictly decreasing on  $(a, b)$ , then there exists  $x_1 < y_1, x_1, y_1 \in (a, b)$  such that

$$(\phi(y_1) - \phi(x_1))(y_1 - x_1) > 0.$$

If for some other pair of points,  $x_2 < y_2$  with  $x_2, y_2 \in (a, b)$ , the above inequality does not hold, then since  $\phi$  is 1-1,

$$(\phi(y_2) - \phi(x_2))(y_2 - x_2) < 0.$$

Let  $x_t \equiv tx_1 + (1 - t)x_2$  and  $y_t \equiv ty_1 + (1 - t)y_2$ . Then  $x_t < y_t$  for all  $t \in [0, 1]$  because

$$tx_1 \leq ty_1 \text{ and } (1 - t)x_2 \leq (1 - t)y_2$$

with strict inequality holding for at least one of these inequalities since not both  $t$  and  $(1 - t)$  can equal zero. Now define

$$h(t) \equiv (\phi(y_t) - \phi(x_t))(y_t - x_t).$$

Since  $h$  is continuous and  $h(0) < 0$ , while  $h(1) > 0$ , there exists  $t \in (0, 1)$  such that  $h(t) = 0$ . Therefore, both  $x_t$  and  $y_t$  are points of  $(a, b)$  and  $\phi(y_t) - \phi(x_t) = 0$  contradicting the assumption that  $\phi$  is one to one. It follows  $\phi$  is either strictly increasing or strictly decreasing on  $(a, b)$ .

This property of being either strictly increasing or strictly decreasing on  $(a, b)$  carries over to  $[a, b]$  by the continuity of  $\phi$ . Suppose  $\phi$  is strictly increasing on  $(a, b)$ , a similar argument holding for  $\phi$  strictly decreasing on  $(a, b)$ . If  $x > a$ , then pick  $y \in (a, x)$  and from the above,  $\phi(y) < \phi(x)$ . Now by continuity of  $\phi$  at  $a$ ,

$$\phi(a) = \lim_{x \rightarrow a+} \phi(x) \leq \phi(y) < \phi(x).$$

Therefore,  $\phi(a) < \phi(x)$  whenever  $x \in (a, b)$ . Similarly  $\phi(b) > \phi(x)$  for all  $x \in (a, b)$ . This proves the lemma.

**Corollary 5.15.4** *Let  $f : (a, b) \rightarrow \mathbb{R}$  be one to one and continuous. Then  $f(a, b)$  is an open interval,  $(c, d)$  and  $f^{-1} : (c, d) \rightarrow (a, b)$  is continuous.*

**Proof:** Since  $f$  is either strictly increasing or strictly decreasing, it follows that  $f(a, b)$  is an open interval,  $(c, d)$ . Assume  $f$  is decreasing. Now let  $x \in (a, b)$ . Why is  $f^{-1}$  continuous at  $f(x)$ ? Since  $f$  is decreasing, if  $f(x) < f(y)$ , then  $y \equiv f^{-1}(f(y)) < x \equiv f^{-1}(f(x))$  and so  $f^{-1}$  is also decreasing. Let  $\varepsilon > 0$  be given. Let  $\varepsilon > \eta > 0$  and  $(x - \eta, x + \eta) \subseteq (a, b)$ . Then  $f(x) \in (f(x + \eta), f(x - \eta))$ . Let  $\delta = \min(f(x) - f(x + \eta), f(x - \eta) - f(x))$ . Then if

$$|f(z) - f(x)| < \delta,$$

it follows

$$z \equiv f^{-1}(f(z)) \in (x - \eta, x + \eta) \subseteq (x - \varepsilon, x + \varepsilon)$$

so

$$|f^{-1}(f(z)) - x| = |f^{-1}(f(z)) - f^{-1}(f(x))| < \varepsilon.$$

This proves the theorem in the case where  $f$  is strictly decreasing. The case where  $f$  is increasing is similar.

# Derivatives

## 6.1 Velocity

Imagine an object which is moving along the real line in the positive direction and that at time  $t > 0$ , the position of the object is  $r(t) = -10 + 30t + t^2$  where distance is measured in kilometers and  $t$  in hours. Thus at  $t = 0$ , the object is at the point  $-10$  kilometers and when  $t = 1$ , the object is at 21 kilometers. The average velocity during this time is the distance traveled divided by the elapsed time. Thus the average velocity would be  $\frac{21 - (-10)}{1} = 31$  kilometers per hour. It came out positive because the object moved in the positive direction along the real line, from  $-10$  to 21. Suppose it was desired to find something which deserves to be referred to as the instantaneous velocity when  $t = 1/2$ ? If the object were a car, it is reasonable to suppose that the magnitude of the average velocity of the object over a very small interval of time would be very close to the number that would appear on the speedometer. For example, if considering the average velocity of the object on the interval  $[\frac{1}{2}, \frac{1}{2} + .0001]$ , this average velocity would be pretty close to the thing which deserves to be called the instantaneous velocity at  $t = \frac{1}{2}$  hours. Thus the velocity at  $t = \frac{1}{2}$  would be close to

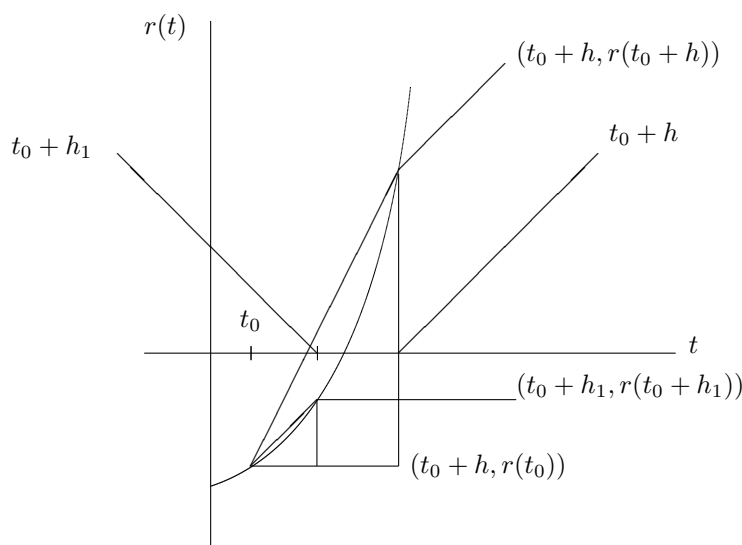
$$\begin{aligned} & (r(.5 + .0001) - r(.5)) / .0001 \\ &= \left( 30(.5 + .01) + (.5 + .0001)^2 - 30(.5) - (.5)^2 \right) / .0001 = 31.0001 \end{aligned}$$

Of course, you would expect to be even closer using a time interval of length .000001 instead of just .0001. In general, consider a time interval of length  $h$  and then define the instantaneous velocity to be the number which all these average velocities get close to as  $h$  gets smaller and smaller. Thus in this case form the average velocity on the interval,  $[\frac{1}{2}, \frac{1}{2} + h]$  to get

$$\left( 30(.5 + h) + (.5 + h)^2 - \left( 30(.5) + (.5)^2 \right) \right) / h = 30 + 2(.5) + h.$$

What number does this average get close to as  $h$  gets smaller and smaller? Clearly it gets close to 31 and for this reason, the velocity at time  $\frac{1}{2}$  is defined as 31. It is positive because the object is moving in the positive direction. If the object were moving in the negative direction, the number would be negative. The notion just described of finding an instantaneous velocity has a geometrical application to finding the slope of a line tangent

to a curve.



In the above picture, you see the slope of the line joining the two points  $(t_0, r(t_0))$  and  $(t_0 + h, r(t_0 + h))$  is given by

$$\frac{r(t_0 + h) - r(t_0)}{h}$$

which equals the average velocity on the time interval,  $[t_0, t_0 + h]$ . You can also see the effect of making  $h$  closer and closer to zero as illustrated by changing  $h$  to the smaller  $h_1$  in the picture. The slope of the resulting line segment appears to get closer and closer to what ought to be considered the slope of the line tangent to the curve at the point  $(t_0, r(t_0))$ .

It is time to make this heuristic material much more precise.

## 6.2 The Derivative

The derivative of a function of one variable is a function given by the following definition.

**Definition 6.2.1** *The derivative of a function,  $f'(x)$ , is defined as the following limit whenever the limit exists. If the limit does not exist, then neither does  $f'(x)$ .*

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \equiv f'(x) \quad (6.1)$$

*The function of  $h$  on the left is called the difference quotient.*

Note that the difference quotient on the left of the equation is a function of  $h$  which is not defined at  $h = 0$ . This is why, in the definition of limit,  $|h| > 0$ . **It is not necessary to have the function defined at the point in order to consider its limit.** The distinction between the limit of a function and its value is very important and must be kept in mind. Also it is clear from setting  $y = x + h$  that

$$f'(x) = \lim_{y \rightarrow x} \frac{f(y) - f(x)}{y - x}. \quad (6.2)$$



**Theorem 6.2.2** *If  $f'(x)$  exists, then  $f$  is continuous at  $x$ .*

**Proof:** Suppose  $\varepsilon > 0$  is given and choose  $\delta_1 > 0$  such that if  $|h| < \delta_1$ ,

$$\left| \frac{f(x+h) - f(x)}{h} - f'(x) \right| < 1.$$

then for such  $h$ , the triangle inequality implies

$$|f(x+h) - f(x)| < |h| + |f'(x)| |h|.$$

Now letting  $\delta < \min\left(\delta_1, \frac{\varepsilon}{1+|f'(x)|}\right)$  it follows if  $|h| < \delta$ , then

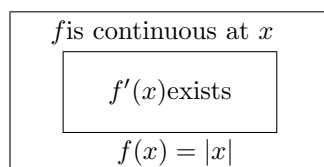
$$|f(x+h) - f(x)| < \varepsilon.$$

Letting  $y = h + x$ , this shows that if  $|y - x| < \delta$ ,

$$|f(y) - f(x)| < \varepsilon$$

which proves  $f$  is continuous at  $x$ .

It is very important to remember that just because  $f$  is continuous, does not mean  $f$  has a derivative. The following picture describes the situation.



As indicated in the above picture the function  $f(x) = |x|$  does not have a derivative at  $x = 0$ . To see this,

$$\lim_{h \rightarrow 0+} \frac{f(h) - f(0)}{h} = \lim_{h \rightarrow 0+} \frac{h}{h} = 1$$

while

$$\lim_{h \rightarrow 0-} \frac{f(h) - f(0)}{h} = \lim_{h \rightarrow 0-} \frac{-h}{h} = -1.$$

Thus the two limits, one from the right and one from the left do not agree as they would have to do if the function had a derivative at  $x = 0$ . See Problem 18 on Page 103. Geometrically, this lack of differentiability is manifested by there being a pointy place in the graph of  $y = |x|$  at  $x = 0$ . In short, the pointy places don't have derivatives.

**Example 6.2.3** *Let  $f(x) = c$  where  $c$  is a constant. Find  $f'(x)$ .*

Set up the difference quotient,

$$\frac{f(x+h) - f(x)}{h} = \frac{c - c}{h} = 0$$

Therefore,

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = \lim_{h \rightarrow 0} 0 = 0$$

**Example 6.2.4** *Let  $f(x) = cx$  where  $c$  is a constant. Find  $f'(x)$ .*

Set up the difference quotient,

$$\frac{f(x+h) - f(x)}{h} = \frac{c(x+h) - cx}{h} = \frac{ch}{h} = c.$$

Therefore,

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = \lim_{h \rightarrow 0} c = c.$$

**Example 6.2.5** Let  $f(x) = \sqrt{x}$  for  $x > 0$ . Find  $f'(x)$ .

Set up the difference quotient,

$$\begin{aligned} \frac{f(x+h) - f(x)}{h} &= \frac{\sqrt{x+h} - \sqrt{x}}{h} = \frac{x+h-x}{h(\sqrt{x+h} + \sqrt{x})} \\ &= \frac{1}{\sqrt{x+h} + \sqrt{x}} \end{aligned}$$

and so

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = \lim_{h \rightarrow 0} \frac{1}{\sqrt{x+h} + \sqrt{x}} = \frac{1}{2\sqrt{x}}.$$

There are rules of derivatives which make finding the derivative very easy.

**Theorem 6.2.6** Let  $a, b \in \mathbb{R}$  and suppose  $f'(t)$  and  $g'(t)$  exist. Then the following formulas are obtained.

$$(af + bg)'(t) = af'(t) + bg'(t). \quad (6.3)$$

$$(fg)'(t) = f'(t)g(t) + f(t)g'(t). \quad (6.4)$$

The formula, (6.4) is referred to as the product rule.

If  $g(t) \neq 0$ ,

$$\left(\frac{f}{g}\right)'(t) = \frac{f'(t)g(t) - g'(t)f(t)}{g^2(t)}. \quad (6.5)$$

Formula (6.5) is referred to as the quotient rule.

If  $f$  is differentiable at  $ct$  where  $c \neq 0$ , Then letting  $g(t) \equiv f(ct)$ ,

$$g'(t) = cf'(ct). \quad (6.6)$$

Written with a slight abuse of notation,

$$(f(ct))' = cf'(ct). \quad (6.7)$$

If  $f$  is differentiable on  $(a, b)$  and if  $g(t) \equiv f(t+c)$ , then  $g$  is differentiable on  $(a-c, b-c)$  and

$$g'(t) = f'(t+c). \quad (6.8)$$

Written with a slight abuse of notation,

$$(f(t+c))' = f'(t+c) \quad (6.9)$$

For  $p$  an integer and  $f'(t)$  exists, let  $g_p(t) \equiv f(t)^p$ . Then

$$(g_p)'(t) = pf(t)^{p-1}f'(t). \quad (6.10)$$

(In the case where  $p < 0$ , assume  $f(t) \neq 0$ .)

Written with a slight abuse of notation, an easy to remember version of (6.10) says

$$(f(t)^p)' = pf(t)^{p-1}f'(t).$$

**Proof:** The first formula is left for you to prove. Consider the second, (6.4).

$$\begin{aligned}\frac{fg(t+h) - fg(t)}{h} &= \frac{f(t+h)g(t+h) - f(t+h)g(t)}{h} + \frac{f(t+h)g(t) - f(t)g(t)}{h} \\ &= f(t+h) \frac{(g(t+h) - g(t))}{h} + \frac{(f(t+h) - f(t))}{h} g(t)\end{aligned}$$

Taking the limit as  $h \rightarrow 0$  and using Theorem 6.2.2 to conclude  $\lim_{h \rightarrow 0} f(t+h) = f(t)$ , it follows from Theorem 5.9.4 that (6.4) follows. Next consider the quotient rule.

$$\begin{aligned}h^{-1} \left( \frac{f}{g}(t+h) - \frac{f}{g}(t) \right) &= \frac{f(t+h)g(t) - g(t+h)f(t)}{hg(t)g(t+h)} \\ &= \frac{f(t+h)(g(t) - g(t+h)) + g(t+h)(f(t+h) - f(t))}{hg(t)g(t+h)} \\ &= \frac{-f(t+h)}{g(t)g(t+h)} \frac{(g(t+h) - g(t))}{h} + \frac{g(t+h)}{g(t)g(t+h)} \frac{(f(t+h) - f(t))}{h}\end{aligned}$$

and from Theorem 5.9.4 on Page 99,

$$\left( \frac{f}{g} \right)'(t) = \frac{g(t)f'(t) - g'(t)f(t)}{g^2(t)}.$$

Now consider Formula (6.6).

$$\begin{aligned}(g(t+h) - g(t))h^{-1} &= h^{-1}(f(ct+ch) - f(ct)) \\ &= c \frac{f(ct+ch) - f(ct)}{ch} \\ &= c \frac{f(ct+h_1) - f(ct)}{h_1}\end{aligned}$$

where  $h_1 = ch$ . Then  $h_1 \rightarrow 0$  if and only if  $h \rightarrow 0$  and so taking the limit as  $h \rightarrow 0$  yields

$$g'(t) = cf'(ct)$$

as claimed. Formulas (6.7) and (6.8) are left as an exercise.

First consider (6.10) in the case where  $p$  equals a nonnegative integer. If  $p = 0$ , (6.10) holds because  $g_0(t) = 1$  and so by Example 6.2.3,

$$g'_0(t) = 0 = 0(f(t))^{-1}f'(t).$$

Next suppose (6.10) holds for  $p$  an integer. Then

$$(g_{p+1}(t)) = f(t)g_p(t)$$

and so by the product rule,

$$\begin{aligned}g'_{p+1}(t) &= f'(t)g_p(t) + f(t)g'_p(t) \\ &= f'(t)(f(t))^p + f(t)(pf(t)^{p-1}f'(t)) \\ &= (p+1)f(t)^p f'(t).\end{aligned}$$

If the formula holds for some integer,  $p$  then it holds for  $-p$ . Here is why.

$$g_{-p}(t) = g_p(t)^{-1}$$

and so

$$\frac{g_{-p}(t+h) - g_{-p}(t)}{h} = \left( \frac{g_p(t) - g_p(t+h)}{h} \right) \left( \frac{1}{g_p(t) g_p(t+h)} \right).$$

Taking the limit as  $h \rightarrow 0$  and using the formula for  $p$ ,

$$\begin{aligned} g'_{-p}(t) &= -p f(t)^{p-1} f'(t) (f(t))^{-2p} \\ &= -p (f(t))^{-p-1} f'(t). \end{aligned}$$

This proves the theorem.

**Example 6.2.7** Let  $p(x) = 3 + 5x + 6x^2 - 7x^3$ . Find  $p'(x)$ .

From the above theorem, and abusing the notation,

$$\begin{aligned} p'(x) &= (3 + 5x + 6x^2 - 7x^3)' \\ &= 3' + (5x)' + (6x^2)' + (-7x^3)' \\ &= 0 + 5 + (6)(2)(x)(x)' + (-7)(3)(x^2)(x)' \\ &= 5 + 12x - 21x^2. \end{aligned}$$

Note the process is to take the exponent and multiply by the coefficient and then make the new exponent one less in each term of the polynomial in order to arrive at the answer. This is the general procedure for differentiating a polynomial as shown in the next example.

**Example 6.2.8** Let  $a_k$  be a number for  $k = 0, 1, \dots, n$  and let  $p(x) = \sum_{k=0}^n a_k x^k$ . Find  $p'(x)$ .

Use Theorem 6.2.6

$$\begin{aligned} p'(x) &= \left( \sum_{k=0}^n a_k x^k \right)' = \sum_{k=0}^n a_k (x^k)' \\ &= \sum_{k=0}^n a_k k x^{k-1} (x)' = \sum_{k=0}^n a_k k x^{k-1} \end{aligned}$$

**Example 6.2.9** Find the derivative of the function  $f(x) = \frac{x^2+1}{x^3}$ .

Use the quotient rule

$$f'(x) = \frac{2x(x^3) - 3x^2(x^2+1)}{x^6} = -\frac{1}{x^4}(x^2+3)$$

**Example 6.2.10** Let  $f(x) = (x^2+1)^4(x^3)$ . Find  $f'(x)$ .

Use the product rule and (6.10). Abusing the notation for the sake of convenience,

$$\begin{aligned} \left( (x^2+1)^4(x^3) \right)' &= \left( (x^2+1)^4 \right)'(x^3) + (x^3)'(x^2+1)^4 \\ &= 4(x^2+1)^3(2x)(x^3) + 3x^2(x^2+1)^4 \\ &= 4x^4(x^2+1)^3 + 3x^2(x^2+1)^4 \end{aligned}$$

**Example 6.2.11** Let  $f(x) = x^3/(x^2+1)^2$ . Find  $f'(x)$ .

Use the quotient rule to obtain

$$f'(x) = \frac{3x^2(x^2+1)^2 - 2(x^2+1)(2x)x^3}{(x^2+1)^4} = \frac{3x^2(x^2+1)^2 - 4x^4(x^2+1)}{(x^2+1)^4}$$

Obviously, one could consider taking the derivative of the derivative and then the derivative of that and so forth. The main thing to consider about this is the notation. The second derivative is denoted with two primes.

**Example 6.2.12** Let  $f(x) = x^3 + 2x^2 + 1$ . Find  $f''(x)$  and  $f'''(x)$ .

To find  $f''(x)$  take the derivative of the derivative. Thus  $f'(x) = 3x^2 + 4x$  and so  $f''(x) = 6x + 4$ . Then  $f'''(x) = 6$ .

When high derivatives are taken, say the 5<sup>th</sup> derivative, it is customary to write  $f^{(5)}(t)$  putting the number of derivatives in parentheses.

## 6.3 Exercises With Answers

1. For  $f(x) = \frac{x^3+6x+2}{x^2+2}$ , find  $f'(x)$ .

Answer:

$$-\frac{-x^4-12+4x}{(x^2+2)^2}$$

2. For  $f(x) = (3x^3 + 6x + 3)(3x^2 + 3x + 9)$ , find  $f'(x)$ .

Answer:

$$(9x^2 + 6)(3x^2 + 3x + 9) + (3x^3 + 6x + 3)(6x + 3)$$

3. For  $f(x) = \sqrt{5x^2 + 1}$ , find  $f'(x)$  from the definition of the derivative.

Answer:

$$\frac{5x}{\sqrt{(5x^2+1)}}$$

4. For  $f(x) = \sqrt[3]{6x^2 + 1}$ , find  $f'(x)$  from the definition of the derivative.

Answer:

$$\frac{4x}{\left(\sqrt[3]{(6x^2+1)}\right)^2}$$

5. For  $f(x) = (-5x + 2)^9$ , find  $f'(x)$  from the definition of the derivative. **Hint:** You might use the formula

$$b^n - a^n = (b - a)(b^{n-1} + b^{n-2}a + \dots + a^{n-2}b + b^{n-1})$$

Answer:

$$-45(-5x + 2)^8$$

6. Let  $f(x) = (x + 5)^2 \sin(1/(x + 5)) + 6(x - 2)(x + 5)$  for  $x \neq -5$  and define  $f(-5) \equiv 0$ . Find  $f'(-5)$  from the definition of the derivative if this is possible. **Hint:** Note that  $|\sin(z)| \leq 1$  for any real value of  $z$ .

Answer:

$$-42$$

## 6.4 Exercises

1. For  $f(x) = -3x^7 + 4x^5 + 2x^3 + x^2 - 5x$ , find  $f^{(3)}(x)$ .
2. For  $f(x) = -x^7 + x^5 + x^3 - 2x$ , find  $f^{(3)}(x)$ .
3. For  $f(x) = -\frac{3x^3 - x - 1}{3x^2 + 1}$ , find  $f'(x)$ .
4. For  $f(x) = (-2x^3 + x + 3)(-2x^2 + 3x - 1)$ , find  $f'(x)$ .
5. For  $f(x) = \sqrt{4x^2 + 1}$ , find  $f'(x)$  from the definition of the derivative.
6. For  $f(x) = \sqrt[3]{3x^2 + 1}$ , find  $f'(x)$  from the definition of the derivative. **Hint:** You might use  

$$b^3 - a^3 = (b^2 + ab + a^2)(b - a) \text{ for } b = \sqrt[3]{3(x+h)^2 + 1} \text{ and } a = \sqrt[3]{3x^2 + 1}.$$
7. For  $f(x) = (-3x + 5)^5$ , find  $f'(x)$  from the definition of the derivative. **Hint:** You might use the formula  

$$b^n - a^n = (b - a)(b^{n-1} + b^{n-2}a + \cdots + a^{n-2}b + b^{n-1})$$
8. Let  $f(x) = (x + 2)^2 \sin(1/(x + 2)) + 6(x - 5)(x + 2)$  for  $x \neq -2$  and define  $f(-2) \equiv 0$ . Find  $f'(-2)$  from the definition of the derivative if this is possible. **Hint:** Note that  $|\sin(z)| \leq 1$  for any real value of  $z$ .
9. Let  $f(x) = (x + 5) \sin(1/(x + 5))$  for  $x \neq -5$  and define  $f(-5) \equiv 0$ . Show  $f'(-5)$  does not exist. **Hint:** Verify that  $\lim_{h \rightarrow 0} \sin(1/h)$  does not exist and then explain why this shows  $f'(-5)$  does not exist.

## 6.5 Local Extrema

When you are on top of a hill, you are at a local maximum although there may be other hills higher than the one on which you are standing. Similarly, when you are at the bottom of a valley, you are at a local minimum even though there may be other valleys deeper than the one you are in. The word, “local” is applied to the situation because if you confine your attention only to points close to your location, you are indeed at either the top or the bottom.

**Definition 6.5.1** Let  $f : D(f) \rightarrow \mathbb{R}$  where here  $D(f)$  is only assumed to be some subset of  $\mathbb{R}$ . Then  $x \in D(f)$  is a local minimum (maximum) if there exists  $\delta > 0$  such that whenever  $y \in (x - \delta, x + \delta) \cap D(f)$ , it follows  $f(y) \geq (\leq) f(x)$ .

Derivatives can be used to locate local maximums and local minimums.

**Theorem 6.5.2** Suppose  $f : (a, b) \rightarrow \mathbb{R}$  and suppose  $x \in (a, b)$  is a local maximum or minimum. Then  $f'(x) = 0$ .

**Proof:** Suppose  $x$  is a local maximum. If  $h > 0$  and is sufficiently small, then  $f(x + h) \leq f(x)$  and so from Theorem 5.9.4 on Page 99,

$$f'(x) = \lim_{h \rightarrow 0^+} \frac{f(x + h) - f(x)}{h} \leq 0.$$

Similarly,

$$f'(x) = \lim_{h \rightarrow 0^-} \frac{f(x + h) - f(x)}{h} \geq 0.$$

The case when  $x$  is local minimum is similar. This proves the theorem.

**Definition 6.5.3** Points where the derivative of a function equals zero are called critical points. It is also customary to refer to points where the derivative of a function does not exist as critical points.

**Example 6.5.4** It is desired to find two positive numbers whose sum equals 16 and whose product is to be as large as possible.

The numbers are  $x$  and  $16 - x$  and  $f(x) = x(16 - x)$  is to be made as large as possible. The value of  $x$  which will do this would be a local maximum so by Theorem 6.5.2 the procedure is to take the derivative of  $f$  and find values of  $x$  where it equals zero. Thus  $16 - 2x = 0$  and the only place this occurs is when  $x = 8$ . Therefore, the two numbers are 8 and 8.

**Example 6.5.5** A farmer wants to fence a rectangular piece of land next to a straight river. What are the dimensions of the largest rectangle if there are exactly 600 meters of fencing available.

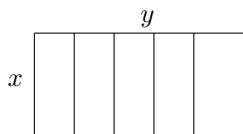
The two sides perpendicular to the river have length  $x$  and the third side has length  $y = (600 - 2x)$ . Thus the function to be maximized is  $f(x) = 2x(600 - x) = 1200x - 2x^2$ . Taking the derivative and setting it equal to zero gives

$$f'(x) = 1200 - 4x = 0$$

and so  $x = 300$ . Therefore, the desired dimensions are  $300 \times 600$ .

**Example 6.5.6** A rectangular playground is to be enclosed by a fence and divided in 5 pieces by 4 fences parallel to one side of the playground. 1704 feet of fencing is used. Find dimensions of the playground which will have the largest total area.

Let  $x$  denote the length of one of these dividing fences and let  $y$  denote the length of the playground as shown in the following picture.



Thus  $6x + 2y = 1704$  so  $y = \frac{1704 - 6x}{2}$  and the function to maximize is

$$f(x) = x \left( \frac{1704 - 6x}{2} \right) = x(852 - 3x).$$

Therefore, to locate the value of  $x$  which will make  $f(x)$  as large as possible, take  $f'(x)$  and set it equal to zero.

$$852 - 6x = 0$$

and so  $x = 142$  feet and  $y = \frac{1704 - 6 \times 142}{2} = 426$  feet.

Revenue is defined to be the amount of money obtained in some transaction. Profit is defined as the revenue minus the costs.

**Example 6.5.7** Francine, the manager of Francine's Fancy Shakes finds that at \$4, demand for her milk shakes is 900 per day. For each \$.30 increase in price, the demand decreases by 50. Find the price and the quantity sold which maximizes revenue.

Let  $x$  be the number of \$.30 increases. Then the total number sold is  $(900 - 50x)$ . The revenue is

$$R(x) = (900 - 50x)(4 + .3x).$$

Then to maximize it,  $R'(x) = 70 - 30x = 0$ . The solution is  $x = 2.33$  and so the optimum price is at \$4.70 and the number sold will be 783.

**Example 6.5.8** *Sam, the owner of Spider Sam's Tarantulas and Creepy Critters finds he can sell 6 tarantulas every day at the regular price of \$30 each. At his last spider celebration sale he reduced the price to \$24 and was able to sell 12 tarantulas every day. He has to pay \$.05 per day to exhibit a tarantula and his fixed costs are \$30 per day, mainly to maintain the thousands of tarantulas he keeps on his tarantula breeding farm. What price should he charge to maximize his profit.*

He assumes the demand for tarantulas is a linear function of price. Thus if  $y$  is the number of tarantulas demanded at price  $x$ , it follows  $y = 36 - x$ . Therefore, the revenue for price  $x$  equals  $R(x) = (36 - x)x$ . Now you have to subtract off the costs to get the profit. Thus

$$P(x) = (36 - x)x - (36 - x)(.05) - 30.$$

It follows the profit is maximized when  $P'(x) = 0$  so  $-2.0x + 36.05 = 0$  which occurs when  $x = \$18.025$ . Thus Sam should charge about \$18 per tarantula.

**Exercise 6.5.9** *Lisa, the owner of Lisa's gags and gadgets sells 500 whoopee cushions per year. It costs \$.25 per year to store a whoopee cushion. To order whoopee cushions it costs \$8 plus \$.90 per cushion. How many times a year and in what lot size should whoopee cushions be ordered to minimize inventory costs?*

Let  $x$  be the times per year an order is sent for a lot size of  $\frac{500}{x}$ . If the demand is constant, it is reasonable to suppose there are about  $\frac{500}{2x}$  whoopee cushions which have to be stored. Thus the cost to store whoopee cushions is  $\frac{125.0}{2x} = .25\left(\frac{500}{2x}\right)$ . Each time an order is made for a lot size of  $\frac{500}{x}$  it costs  $8 + \frac{450.0}{x} = 8 + .9\left(\frac{500}{x}\right)$  and this is done  $x$  times a year. Therefore, the total inventory cost is  $x\left(8 + \frac{450.0}{x}\right) + \frac{125.0}{2x} = C(x)$ . The problem is to minimize  $C(x) = x\left(8 + \frac{450.0}{x}\right) + \frac{125.0}{2x}$ . Taking the derivative yields

$$C'(x) = \frac{16x^2 - 125}{2x^2}$$

and so the value of  $x$  which will minimize  $C(x)$  is  $\frac{5}{4}\sqrt{5} = 2.79 \dots$  and the lot size is  $\frac{500}{\frac{5}{4}\sqrt{5}} = 178.88$ . Of course you would round these numbers off. Order 179 whoopee cushions 3 times a year.

## 6.6 Exercises With Answers

- Find the  $x$  values of the critical points of the function  $f(x) = x^2 - x^3$ .

Answer:

$$\frac{2}{3}, 0$$

- Find the  $x$  values of the critical points of the function  $f(x) = \sqrt{3x^2 - 6x + 6}$ .

Answer:

$$1$$



3. Find the extreme points of the function,  $f(x) = 2x^2 - 20x + 55$  and tell whether the extreme point is a maximum or a minimum.

Answer:

The extremum is at  $x = 5$ . It is a maximum.

4. Find the extreme points of the function,  $f(x) = x + \frac{9}{x}$  and tell whether the extreme point is a local maximum or a local minimum or neither.

Answer:

The extrema are at  $x = \pm 3$ . The one at 3 is a local minimum and the one at  $-3$  is a local maximum.

5. A rectangular pasture is to be fenced off beside a river with no need of fencing along the river. If there is 900 yards of fencing material, what are the dimensions of the largest possible pasture that can be enclosed?

Answer:

$225 \times 45$

6. A piece of property is to be fenced on the front and two sides. Fencing for the sides costs \$3.50 per foot and fencing for the front costs \$5.60 per foot. What are the dimensions of the largest such rectangular lot if the available money is \$840.0?

Answer:

$60 \times 75$ .

7. In a particular apartment complex of 120 units, it is found that all units remain occupied when the rent is \$300 per month. For each \$30 increase in the rent, one unit becomes vacant, on the average. Occupied units require \$60 per month for maintenance, while vacant units require none. Fixed costs for the buildings are \$30 000 per month. What rent should be charged for maximum profit and what is the maximum profit?

Answer:

Need to maximize  $f(x) = (300 + 30x)(120 - x) - 37\,200 + 60x$  for  $x \in [0, 120]$  where  $x$  is the number of \$30 increases in rent.

\$92 880 when the rent is \$1980.

8. A picture is 5 feet high and the eye level of an observer is 2 feet below the bottom edge of the picture. How far from the picture should the observer stand if he wants to maximize the angle subtended by the picture?

Answer:

Let the angle subtended by the picture be  $\theta$  and let  $\alpha$  denote the angle between a horizontal line from the observer's eye to the wall and the line between the observer's eye and the base of the picture. Then letting  $x$  denote the distance between the wall and the observer's eye,  $7 = x \tan(\theta + \alpha) = x \left( \frac{\tan \theta + \tan \alpha}{1 - \tan \theta \tan \alpha} \right) = x \left( \frac{\tan \theta + \frac{2}{x}}{1 - \tan \theta \left( \frac{2}{x} \right)} \right)$ . The problem is equivalent to maximizing  $\tan \theta$  so denote this by  $z$  and solve for it. Thus  $x \left( \frac{z + \frac{2}{x}}{1 - z \left( \frac{2}{x} \right)} \right) = 7$  and so  $z = 5 \frac{x}{14 + x^2}$ . It follows  $\frac{dz}{dx} = 5 \frac{14 - x^2}{(14 + x^2)^2}$  and setting this equal to zero,  $x = \sqrt{14}$ .

9. Find the point on the curve,  $y = \sqrt{81 - 6x}$  which is closest to  $(0, 0)$ .

Answer:

$$(3, \sqrt{63})$$

10. A street is 200 feet long and there are two lights located at the ends of the street. One of the lights is  $\frac{27}{8}$  times as bright as the other. Assuming the brightness of light from one of these street lights is proportional to the brightness of the light and the reciprocal of the square of the distance from the light, locate the darkest point on the street.

Answer:

80 feet from one light and 120 feet from the other.

11. Two cities are located on the same side of a straight river. One city is at a distance of 3 miles from the river and the other city is at a distance of 8 miles from the river. The distance between the two points on the river which are closest to the respective cities is 40 miles. Find the location of a pumping station which is to pump water to the two cities which will minimize the length of pipe used.

Answer:

$\frac{120}{11}$  miles from the point on the river closest to the city which is at a distance of 3 miles from the river.

## 6.7 Exercises

- If  $f'(x) = 0$ , is it necessary that  $x$  is either a local minimum or local maximum? **Hint:** Consider  $f(x) = x^3$ .
- Two positive numbers add to 32. Find the numbers if their product is to be as large as possible.
- The product of two positive numbers equals 16. Find the numbers if their sum is to be as small as possible.
- The product of two positive numbers equals 16. Find the numbers if twice the first plus three times the second is to be as small as possible.
- Theodore, the owner of Theodore's tarantulas finds he can sell 6 tarantulas at the regular price of \$20 each. At his last spider celebration day sale he reduced the price to \$14 and was able to sell 14 tarantulas. He has to pay \$.05 per day to maintain a tarantula and his fixed costs are \$30 per day. What price should he charge to maximize his profit.
- Lisa, the owner of Lisa's gags and gadgets sells 500 whoopee cushions per year. It costs \$.25 per year to store a whoopee cushion. To order whoopee cushions it costs \$2 plus \$.25 per cushion. How many times a year and in what lot size should whoopee cushions be ordered to minimize inventory costs?
- A continuous function,  $f$  defined on  $[a, b]$  is to be maximized. It was shown above in Theorem 6.5.2 that if the maximum value of  $f$  occurs at  $x \in (a, b)$ , and if  $f$  is differentiable there, then  $f'(x) = 0$ . However, this theorem does not say anything about the case where the maximum of  $f$  occurs at either  $a$  or  $b$ . Describe how to find the point of  $[a, b]$  where  $f$  achieves its maximum. Does  $f$  have a maximum? Explain.

8. Find the maximum and minimum values and the values of  $x$  where these are achieved for the function,  $f(x) = x + \sqrt{25 - x^2}$ .
9. A piece of wire of length  $L$  is to be cut in two pieces. One piece is bent into the shape of an equilateral triangle and the other piece is bent to form a square. How should the wire be cut to maximize the sum of the areas of the two shapes? How should the wire be bent to minimize the sum of the areas of the two shapes? **Hint:** Be sure to consider the case where all the wire is devoted to one of the shapes separately. This is a possible solution even though the derivative is not zero there.
10. A cylindrical can is to be constructed of material which costs 3 cents per square inch for the top and bottom and only 2 cents per square inch for the sides. The can needs to hold  $90\pi$  cubic inches. Find the dimensions of the cheapest can. **Hint:** The volume of a cylinder is  $\pi r^2 h$  where  $r$  is the radius of the base and  $h$  is the height. The area of the cylinder is  $2\pi r^2 + 2\pi r h$ .
11. A rectangular sheet of tin has dimensions 10 cm. by 20 cm. It is desired to make a topless box by cutting out squares from each corner of the rectangular sheet and then folding the rectangular tabs which remain. Find the volume of the largest box which can be made in this way.
12. Let  $f(x) = \frac{1}{3}x^3 - x^2 - 8x$  on the interval  $[-1, 10]$ . Find the point of  $[-1, 10]$  at which  $f$  achieves its minimum.
13. A rectangular garden 200 square feet in area is to be fenced off against rabbits. Find the least possible length of fencing if one side of the garden is already protected by a barn.
14. A feed lot is to be enclosed by a fence and divided in 5 pieces by 4 fences parallel to one side. 1272 feet of fencing is used. Find dimensions of the feed lot which will have the largest total area.
15. Find the dimensions of the largest rectangle that can be inscribed in a semicircle of radius  $r$  where  $r = 8$ .
16. A smuggler wants to fit a small cylindrical vial inside a hollow rubber ball with a eight inch diameter. Find the volume of the largest vial that can fit inside the ball. The volume of a cylinder equals  $\pi r^2 h$  where  $h$  is the height and  $r$  is the radius.
17. A function,  $f$ , is said to be odd if  $f(-x) = -f(x)$  and a function is said to be even if  $f(-x) = f(x)$ . Show that if  $f'$  is even, then  $f$  is odd and if  $f'$  is odd, then  $f$  is even. Sketch the graph of a typical odd function and a typical even function.
18. Recall  $\sin$  is an odd function and  $\cos$  is an even function. Determine whether each of the trig functions is odd, even or neither.
19. Find the  $x$  values of the critical points of the function  $f(x) = 3x^2 - 5x^3$ .
20. Find the  $x$  values of the critical points of the function  $f(x) = \sqrt{3x^2 - 6x + 8}$ .
21. Find the extreme points of the function,  $f(x) = x + \frac{25}{x}$  and tell whether the extreme point is a local maximum or a local minimum or neither.
22. A piece of property is to be fenced on the front and two sides. Fencing for the sides costs \$3.50 per foot and fencing for the front costs \$5.60 per foot. What are the dimensions of the largest such rectangular lot if the available money is \$1400?

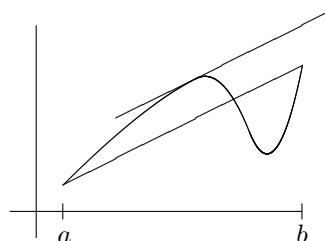
23. In a particular apartment complex of 200 units, it is found that all units remain occupied when the rent is \$400 per month. For each \$40 increase in the rent, one unit becomes vacant, on the average. Occupied units require \$80 per month for maintenance, while vacant units require none. Fixed costs for the buildings are \$20 000 per month. What rent should be charged for maximum profit and what is the maximum profit?
24. A picture is 9 feet high and the eye level of an observer is 2 feet below the bottom edge of the picture. How far from the picture should the observer stand if he wants to maximize the angle subtended by the picture?
25. Find the point on the curve,  $y = \sqrt{25 - 2x}$  which is closest to  $(0, 0)$ .
26. A street is 200 feet long and there are two lights located at the ends of the street. One of the lights is  $\frac{1}{8}$  times as bright as the other. Assuming the brightness of light from one of these street lights is proportional to the brightness of the light and the reciprocal of the square of the distance from the light, locate the darkest point on the street.
27. Two cities are located on the same side of a straight river. One city is at a distance of 3 miles from the river and the other city is at a distance of 8 miles from the river. The distance between the two points on the river which are closest to the respective cities is 40 miles. Find the location of a pumping station which is to pump water to the two cities which will minimize the length of pipe used.

## 6.8 Mean Value Theorem

The mean value theorem is one of the most important theorems about the derivative. The best versions of many other theorems depend on this fundamental result. The mean value theorem says that under suitable conditions, there exists a point in  $(a, b)$ ,  $x$ , such that  $f'(x)$  equals the slope of the secant line,

$$\frac{f(b) - f(a)}{b - a}.$$

The following picture is descriptive of this situation.



This theorem is an existence theorem and like the other existence theorems in analysis, it depends on the completeness axiom. The following is known as Rolle's<sup>1</sup> theorem.

**Theorem 6.8.1** Suppose  $f : [a, b] \rightarrow \mathbb{R}$  is continuous,

$$f(a) = f(b),$$

---

<sup>1</sup>Rolle is remembered for Rolle's theorem and not for anything else he did. Ironically, he did not like calculus.

and

$$f : (a, b) \rightarrow \mathbb{R}$$

has a derivative at every point of  $(a, b)$ . Then there exists  $x \in (a, b)$  such that  $f'(x) = 0$ .

**Proof:** Suppose first that  $f(x) = f(a)$  for all  $x \in [a, b]$ . Then any  $x \in (a, b)$  is a point such that  $f'(x) = 0$ . If  $f$  is not constant, either there exists  $y \in (a, b)$  such that  $f(y) > f(a)$  or there exists  $y \in (a, b)$  such that  $f(y) < f(a)$ . In the first case, the maximum of  $f$  is achieved at some  $x \in (a, b)$  and in the second case, the minimum of  $f$  is achieved at some  $x \in (a, b)$ . Either way, Theorem 6.5.2 on Page 126 implies  $f'(x) = 0$ . This proves Rolle's theorem.

The next theorem is known as the Cauchy mean value theorem.

**Theorem 6.8.2** Suppose  $f, g$  are continuous on  $[a, b]$  and differentiable on  $(a, b)$ . Then there exists  $x \in (a, b)$  such that

$$f'(x)(g(b) - g(a)) = g'(x)(f(b) - f(a)).$$

**Proof:** Let

$$h(x) \equiv f(x)(g(b) - g(a)) - g(x)(f(b) - f(a)).$$

Then letting  $x = a$  and then letting  $x = b$ , a short computation shows  $h(a) = h(b)$ . Also,  $h$  is continuous on  $[a, b]$  and differentiable on  $(a, b)$ . Therefore Rolle's theorem applies and there exists  $x \in (a, b)$  such that

$$h'(x) = f'(x)(g(b) - g(a)) - g'(x)(f(b) - f(a)) = 0.$$

This proves the theorem.

The usual mean value theorem, sometimes called the Lagrange mean value theorem, illustrated by the above picture is obtained by letting  $g(x) = x$ .

**Corollary 6.8.3** Let  $f$  be continuous on  $[a, b]$  and differentiable on  $(a, b)$ . Then there exists  $x \in (a, b)$  such that  $f(b) - f(a) = f'(x)(b - a)$ .

**Corollary 6.8.4** Suppose  $f'(x) = 0$  for all  $x \in (a, b)$  where  $a \geq -\infty$  and  $b \leq \infty$ . Then  $f(x) = f(y)$  for all  $x, y \in (a, b)$ . Thus  $f$  is a constant.

**Proof:** If this is not true, there exists  $x_1$  and  $x_2$  such that  $f(x_1) \neq f(x_2)$ . Then by the mean value theorem,

$$0 \neq \frac{f(x_1) - f(x_2)}{x_1 - x_2} = f'(z)$$

for some  $z$  between  $x_1$  and  $x_2$ . This contradicts the hypothesis that  $f'(x) = 0$  for all  $x$ . This proves the theorem.

**Corollary 6.8.5** Suppose  $f'(x) > 0$  for all  $x \in (a, b)$  where  $a \geq -\infty$  and  $b \leq \infty$ . Then  $f$  is strictly increasing on  $(a, b)$ . That is, if  $x < y$ , then  $f(x) < f(y)$ . If  $f'(x) \geq 0$ , then  $f$  is increasing in the sense that whenever  $x < y$  it follows that  $f(x) \leq f(y)$ .

**Proof:** Let  $x < y$ . Then by the mean value theorem, there exists  $z \in (x, y)$  such that

$$0 < f'(z) = \frac{f(y) - f(x)}{y - x}.$$

Since  $y > x$ , it follows  $f(y) > f(x)$  as claimed. Replacing  $<$  by  $\leq$  in the above equation and repeating the argument gives the second claim.

**Corollary 6.8.6** Suppose  $f'(x) < 0$  for all  $x \in (a, b)$  where  $a \geq -\infty$  and  $b \leq \infty$ . Then  $f$  is strictly decreasing on  $(a, b)$ . That is, if  $x < y$ , then  $f(x) > f(y)$ . If  $f'(x) \leq 0$ , then  $f$  is decreasing in the sense that for  $x < y$ , it follows that  $f(x) \geq f(y)$ .

**Proof:** Let  $x < y$ . Then by the mean value theorem, there exists  $z \in (x, y)$  such that

$$0 > f'(z) = \frac{f(y) - f(x)}{y - x}.$$

Since  $y > x$ , it follows  $f(y) < f(x)$  as claimed. The second claim is similar except instead of a strict inequality in the above formula, you put  $\geq$ .

## 6.9 Exercises

1. Sally drives her Saturn over the 110 mile toll road in exactly 1.3 hours. The speed limit on this toll road is 70 miles per hour and the fine for speeding is 10 dollars per mile per hour over the speed limit. How much should Sally pay?
2. Two cars are careening down a freeway weaving in and out of traffic. Car A passes car B and then car B passes car A as the driver makes obscene gestures. This infuriates the driver of car A who passes car B while firing his handgun at the driver of car B. Show there are at least two times when both cars have the same speed. Then show there exists at least one time when they have the same acceleration. The acceleration is the derivative of the velocity.
3. Show the cubic function,  $f(x) = 5x^3 + 7x - 18$  has only one real zero.
4. Suppose  $f(x) = x^7 + |x| + x - 12$ . How many solutions are there to the equation,  $f(x) = 0$ ?
5. Let  $f(x) = |x - 7| + (x - 7)^2 - 2$  on the interval  $[6, 8]$ . Then  $f(6) = 0 = f(8)$ . Does it follow from Rolle's theorem that there exists  $c \in (6, 8)$  such that  $f'(c) = 0$ ? Explain your answer.
6. Suppose  $f$  and  $g$  are differentiable functions defined on  $\mathbb{R}$ . Suppose also that it is known that  $|f'(x)| > |g'(x)|$  for all  $x$  and that  $|f'(t)| > 0$  for all  $t$ . Show that whenever  $x \neq y$ , it follows  $|f(x) - f(y)| > |g(x) - g(y)|$ . **Hint:** Use the Cauchy mean value theorem, Theorem 6.8.2.
7. Show that, like continuous functions, functions which are derivatives have the intermediate value property. This means that if  $f'(a) < 0 < f'(b)$  then there exists  $x \in (a, b)$  such that  $f'(x) = 0$ . **Hint:** Argue the minimum value of  $f$  occurs at an interior point of  $[a, b]$ .
8. Consider the function

$$f(x) \equiv \begin{cases} 1 & \text{if } x \geq 0 \\ -1 & \text{if } x < 0 \end{cases}.$$

Is it possible that this function could be the derivative of some function? Why?

## 6.10 Curve Sketching

The theorems and corollaries given above can be used to aid in sketching the graphs of functions. The second derivative will also help in determining the shape of the function.

**Definition 6.10.1** A differentiable function,  $f$ , defined on an interval,  $(a, b)$ , is concave up if  $f'$  is an increasing function. A differentiable function defined on an interval,  $(a, b)$ , is concave down if  $f'$  is a decreasing function.

From the geometric description of the derivative as the slope of a tangent line to the graph of the function, to say derivative is an increasing function means that as you move from left to right, the slopes of the lines tangent to the graph of  $f$  become larger. Thus the graph of the function is bent up in the shape of a smile. It may also help to think of it as a cave when you view it from above, hence the term concave up. If the derivative is decreasing, it follows that as you move from left to right the slopes of the lines tangent to the graph of  $f$  become smaller. Thus the graph of the function is bent down in the form of a frown. It is concave down because it is like a cave when viewed from beneath. The following theorem will give a convenient criterion in terms of the second derivative for finding whether a function is concave up or concave down. The term, concavity, is used to refer to this property. Thus you determine the concavity of a function when you find whether it is concave up or concave down.

**Theorem 6.10.2** Suppose  $f''(x) > 0$  for  $x \in (a, b)$ . Then  $f$  is concave up on  $(a, b)$ . Suppose  $f''(x) < 0$  on  $(a, b)$ . Then  $f$  is concave down.

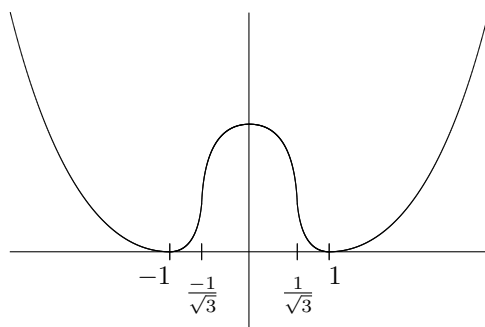
**Proof:** This follows immediately from Corollaries 6.8.6 and 6.8.5 applied to the first derivative. The following picture may help in remembering this.



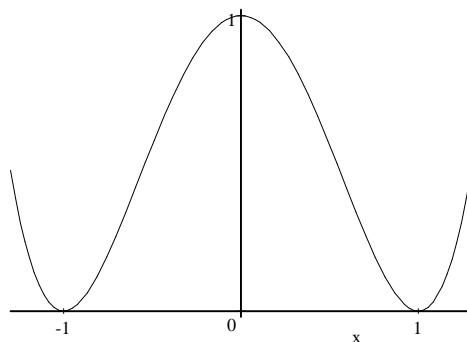
In this picture, the plus signs and the smile on the left correspond to the second derivative being positive. The smile gives the way in which the graph of the function is bent. In the second face, the minus signs correspond to the second derivative being negative. The frown gives the way in which the graph of the function is bent.

**Example 6.10.3** Sketch the graph of the function,  $f(x) = (x^2 - 1)^2 = x^4 - 2x^2 + 1$

Take the derivative of this function,  $f'(x) = 4x^3 - 4x = 4x(x - 1)(x + 1)$  which equals zero at  $-1, 0$ , and  $1$ . It is positive on  $(-1, 0)$ , and  $(1, \infty)$  and negative on  $(0, 1)$  and  $(-\infty, -1)$ . Therefore,  $x = 0$  corresponds to a local maximum and  $x = -1$  and  $x = 1$  correspond to local minimums. The second derivative is  $f''(x) = 12x^2 - 4$  and this equals zero only at the points  $-1/\sqrt{3}$  and  $1/\sqrt{3}$ . The second derivative is positive on the intervals  $(1/\sqrt{3}, \infty)$  and  $(-\infty, -1/\sqrt{3})$  so the function,  $f$  is smiling on these intervals. The second derivative is negative on the interval  $(-1/\sqrt{3}, 1/\sqrt{3})$  and so the original function is frowning on this interval. This describes in words the qualitative shape of the function. It only remains to draw a picture which incorporates this description. The following is such a sketch. It is not intended to be an accurate drawing made to scale, only to be a qualitative picture of what was just determined.



A better graph of this function is the following, done by a computer algebra system. However the computer worked a lot harder.



In general, if you are interested in getting a nice graph of a function, you should use a computer algebra system. An effective way to accomplish your graphing is to go to the help menu and copy and paste an example from this menu changing it as needed. Both mathematica and Maple have good help menus. Keep in mind there are certain conventions which must be followed. For example to write  $x$  raised to the second power you enter  $x^2$ . In Maple, you also need to place an asterisk between quantities which are multiplied since otherwise it will not know you are multiplying and won't work. There are also easy to use versions of Maple available which involve essentially pointing and clicking. You won't learn any calculus from playing with a computer algebra system but you might have a lot of fun.

## 6.11 Exercises

1. Sketch the graph of the function,  $f(x) = x^3 - 3x + 1$  showing the intervals on which the function is concave up and down and identifying the intervals on which the function is increasing.
2. Find intervals on which the function,  $f(x) = \sqrt{1 - x^2}$  is increasing and intervals on which it is concave up and concave down. Sketch a graph of the function.
3. Sketch the graphs of  $y = x^4$ ,  $y = x^3$ , and  $y = -x^4$ . What do these graphs tell you about the case when the second derivative equals zero?
4. Sketch the graph of  $f(x) = 1/(1 + x^2)$  showing the intervals on which the function is increasing or decreasing and the intervals on which the graph is concave up and



concave down.

5. Sketch the graph of  $f(x) = x/(1+x^2)$  showing the intervals on which the function is increasing or decreasing and the intervals on which the graph is concave up and concave down.
6. Inflection points are points where the graph of a function changes from being concave up to concave down or from being concave down to being concave up. Show that inflection points can be identified by looking at those points where the second derivative equals zero but that not every point where the second derivative equals zero is an inflection point. **Hint:** For the last part consider  $y = x^3$  and  $y = x^4$ .
7. Find all inflection points for the function,  $f(x) = x^2/(1+x^2)$ .



# Some Important Special Functions

## 7.1 The Circular Functions

The Trigonometric functions are also called the circular functions. Thus this section will be on the functions,  $\cos$ ,  $\sin$ ,  $\tan$ ,  $\sec$ ,  $\csc$ , and  $\cot$ . The first thing to do is to give an important lemma. There are several approaches to this lemma. To see it done in terms of areas of a circular sector, see Apostol, [2], or almost any other calculus book. However, the book by Apostol has no loose ends in the presentation unlike most other books which use this approach. The proof given here is a modification of that found in Tierney, [17] and Rose, [13] and is based on arc length.

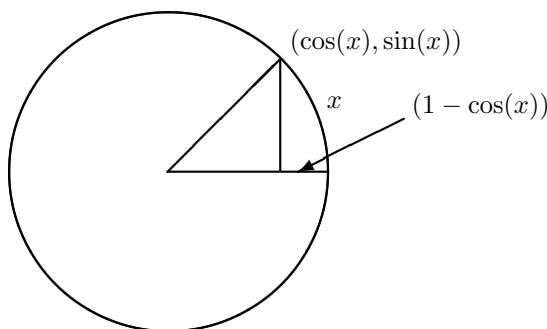
**Lemma 7.1.1** *The following limits hold.*

$$\lim_{x \rightarrow 0} \frac{\sin x}{x} = 1 \quad (7.1)$$

$$\lim_{x \rightarrow 0} \frac{1 - \cos x}{x} = 0 \quad (7.2)$$

**Proof:** First consider (7.1). In the following picture, it follows from Corollary 3.5.4 on Page 59 that for small positive  $x$ ,

$$\sin x + (1 - \cos x) \geq x \geq \sin x. \quad (7.3)$$



Now divide by  $\sin x$  to get

$$1 + \frac{1 - \cos x}{|\sin x|} = 1 + \frac{1 - \cos x}{\sin x} \geq \frac{x}{\sin x} \geq 1.$$

For small negative values of  $x$ , it is also true that

$$1 + \frac{1 - \cos x}{|\sin x|} \geq \frac{x}{\sin x} \geq 1.$$

(Why?) From the trig. identities, it follows that for all small values of  $x$ ,

$$1 + \frac{\sin^2 x}{|\sin x| (1 + \cos x)} = 1 + \frac{|\sin x|}{(1 + \cos x)} \geq \frac{x}{\sin x} \geq 1$$

and so from the squeezing theorem, Theorem 5.9.5 on Page 100,

$$\lim_{x \rightarrow 0} \frac{x}{\sin x} = 1$$

and consequently, from the limit theorems,

$$\lim_{x \rightarrow 0} \frac{\sin x}{x} = \lim_{x \rightarrow 0} \frac{1}{\left(\frac{x}{\sin x}\right)} = 1.$$

Finally,

$$\frac{1 - \cos x}{x} = \frac{1 - \cos^2 x}{x(1 + \cos x)} = \sin x \frac{\sin x}{x} \frac{1}{1 + \cos x}.$$

Therefore, from Theorem 5.5.1 on Page 92 which says  $\lim_{x \rightarrow 0} \sin(x) = 0$ , and the limit theorems,

$$\lim_{x \rightarrow 0} \frac{1 - \cos x}{x} = 0.$$

This proves the Lemma.

With this, it is easy to find the derivative of  $\sin$ . Using Lemma 7.1.1,

$$\begin{aligned} \lim_{h \rightarrow 0} \frac{\sin(x+h) - \sin x}{h} &= \lim_{h \rightarrow 0} \frac{\sin(x) \cos(h) + \cos(x) \sin(h) - \sin x}{h} \\ &= \lim_{h \rightarrow 0} \frac{(\sin x)(\cos(h) - 1)}{h} + \cos x \frac{\sin(h)}{h} \\ &= \cos x. \end{aligned}$$

The derivative of  $\cos$  can be found the same way. Alternatively,

$$\cos(x) = \sin(x + \pi/2)$$

and so

$$\begin{aligned} \cos'(x) &= \sin'(x + \pi/2) \\ &= \cos(x + \pi/2) \\ &= \cos x \cos(\pi/2) - \sin x \sin(\pi/2) \\ &= -\sin x. \end{aligned}$$

The following theorem is now obvious and the proofs of the remaining parts are left for you.

**Theorem 7.1.2** *The derivatives of the trig. functions are as follows.*

$$\begin{aligned} \sin'(x) &= \cos x \\ \cos'(x) &= -\sin x \\ \tan'(x) &= \sec^2(x) \\ \cot'(x) &= -\csc^2(x) \\ \sec'(x) &= \sec x \tan x \\ \csc'(x) &= -\csc x \cot x \end{aligned}$$

Here are some examples of extremum problems which involve the use of the trig. functions.

**Example 7.1.3** *Two hallways intersect at a right angle. One is 5 feet wide and the other is 2 feet wide. What is the length of the longest thin rod which can be carried horizontally from one hallway to the other?*

You must minimize the length of the rod which touches the inside corner of the two halls and extends to the outside walls. Letting  $\theta$  be the angle between this rod and the outside wall for the hall having width 2, minimize

$$f(\theta) = \overbrace{2 \csc \theta + 5 \sec \theta}^{\text{length of rod}}.$$

Therefore, using the rules of differentiation,

$$f'(\theta) = \frac{2 \cos^3 \theta - 5 \sin \theta + 5 \sin \theta \cos^2 \theta}{(\cos^2 \theta)(-1 + \cos^2 \theta)} = 0$$

should be solved to get the angle where this length is as small as possible. Thus

$$2 \cos^3 \theta - 5 \sin \theta + 5 \sin \theta \cos^2 \theta = 0.$$

and  $2 \cos^3 \theta - 5 \sin^3 \theta = 0$  and so  $\tan \theta = \frac{1}{5} \sqrt[3]{2} (\sqrt[3]{5})^2$ . Drawing a triangle, you see that at this value of  $\theta$ , you have  $\sec \theta = \frac{\sqrt{(\sqrt[3]{5})^2 + (\sqrt[3]{2})^2}}{\sqrt[3]{5}}$  and  $\csc \theta = \frac{\sqrt{(\sqrt[3]{5})^2 + (\sqrt[3]{2})^2}}{\sqrt[3]{2}}$ . Therefore, the minimum is obtained by substituting these values in to the equation for  $f(\theta)$  yielding  $\left( \sqrt{((\sqrt[3]{5})^2 + (\sqrt[3]{2})^2)} \right)^3$ .

**Example 7.1.4** *A fence 9 feet high is 2 feet from a building. What is the length of the shortest ladder which will lean against the top of the fence and touch the building?*

Let  $\theta$  be the angle of the ladder with the ground. Then the length of this ladder making this angle with the ground and leaning on the top of the fence while touching the building is

$$f(\theta) = \frac{9}{\sin \theta} + \frac{2}{\cos \theta}.$$

Then the final answer is  $\left( \sqrt{3\sqrt[3]{3} + (\sqrt[3]{2})^2} \right)^3$ . The details are similar to the problem of the two hallways.

## 7.2 Exercises

1. Prove all parts of Theorem 7.1.2.
2. Prove  $\tan'(x) = 1 + \tan^2(x)$ .
3. Find and prove a formula for the derivative of  $\sin^m(x)$  for  $m$  an integer.
4. Find the derivative of the function,  $\sin^6(5x)$ .
5. Find the derivative of the function,  $\tan^7(4x)$ .

6. Find the derivative of the function,  $\frac{\sec^3(2x)}{\tan^3(3x)}$ .
7. Find all intervals where  $\sin(2x)$  is concave down.
8. Find the intervals where  $\cos(3x)$  is increasing.
9. Two hallways intersect at a right angle. One is 3 feet wide and the other is 4 feet wide. What is the length of the longest thin rod which can be carried horizontally from one hallway to the other?
10. A fence 5 feet high is 2 feet from a building. What is the length of the shortest ladder which will lean against the top of the fence and touch the building?
11. Suppose  $f(x) = A \cos \omega x + B \sin \omega x$ . Show there exists an angle,  $\phi$  such that  $f(x) = \sqrt{A^2 + B^2} \sin(\omega x + \phi)$ . The number,  $\sqrt{A^2 + B^2}$  gives the “amplitude” and  $\phi$  is called the “phase shift” while  $\omega$  is called the “frequency”. This is very important because it allows us to understand what is going on. The amplitude gives the height of the periodic function,  $f$ . **Hint:** Remember a point on the unit circle determines an angle. Write  $f(x)$  in the form

$$\sqrt{A^2 + B^2} \left( \frac{A}{\sqrt{A^2 + B^2}} \cos \omega x + \frac{B}{\sqrt{A^2 + B^2}} \sin \omega x \right)$$

and note that  $\left( \frac{B}{\sqrt{A^2 + B^2}}, \frac{A}{\sqrt{A^2 + B^2}} \right)$  is a point on the unit circle.

12. Repeat Problem 11 but this time show  $f(x) = \sqrt{A^2 + B^2} \cos(\omega x + \phi)$ . How could you find  $\phi$ ?

## 7.3 The Exponential And Log Functions

### 7.3.1 The Rules Of Exponents

As mentioned earlier,  $b^m$  means to multiply  $b$  by itself  $m$  times assuming  $m$  is a positive integer.  $b^0 \equiv 1$  provided  $b \neq 0$ . In the case where  $b = 0$  the symbol is undefined. If  $m < 0$ ,  $b^m$  is defined as  $\frac{1}{b^{-m}}$ . Then the following algebraic properties are obtained. Be sure you understand these properties for  $x$  and  $y$  integers.

$$b^{x+y} = b^x b^y, (ab)^x = a^x b^x \quad (7.4)$$

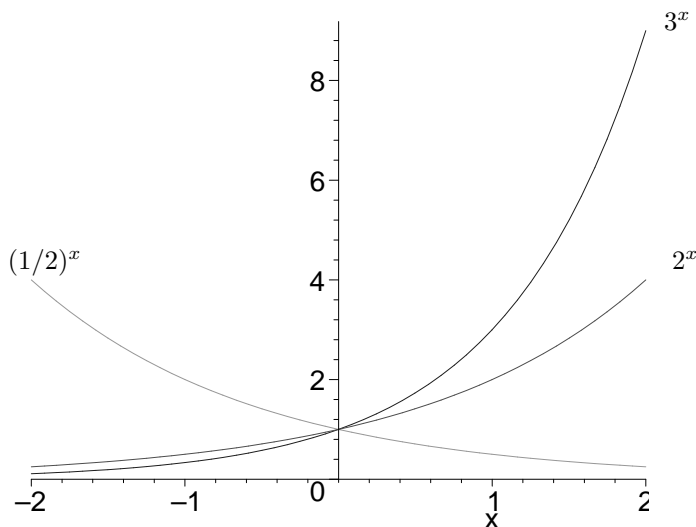
$$b^{xy} = (b^x)^y, b^{-1} = \frac{1}{b} \quad (7.5)$$

These properties are called the rules of exponents.

When  $x$  and  $y$  are not integers, the meaning of  $b^x$  is no longer clear. For example, suppose  $b = -1$  and  $x = 1/2$ . What exactly is meant by  $(-1)^{1/2}$ ? Even in the case where  $b > 0$  there are difficulties. If  $x$  is a rational number,  $m/n$  and  $b > 0$  the symbol  $b^{m/n}$  means  $\sqrt[n]{b^m}$ . That is its definition and it is a useful exercise for you to verify (7.4) and (7.5) hold with this definition. There are no mathematical questions about the existence of this number. To see this, consider Problem 6 on Page 97. The problem is not one of theory but of practicality. Could you use this definition to find  $2^{\frac{1234567812345}{1234567812344}}$ ? Consider what you would do. First find the number  $2^{1234567812345}$  and then  $\dots$ ? Can you find this number? It is just too big. However, a calculator can find  $2^{\frac{1234567812345}{1234567812344}}$ . It yields  $2^{\frac{1234567812345}{1234567812344}} = 2.000\,000\,000\,001\,123$  as an approximate answer. Clearly something else must be going on. To make matters even worse, what would you do with  $2^{\sqrt{2}}$ ? As mentioned earlier,  $\sqrt{2}$  is irrational and so cannot be written as the quotient of two integers. These are serious difficulties and must be dealt with.

### 7.3.2 The Exponential Functions, A Wild Assumption

Using your calculator or a computer you can obtain graphs of the functions,  $y = b^x$  for various choices of  $b$ . The following picture gives a few of these graphs.



These graphs suggest that if  $b < 1$  the function,  $y = b^x$  is decreasing while if  $b > 1$ , the function is increasing but just how was the calculator or computer able to draw those graphs? Also, do the laws of exponents continue to hold for all real values of  $x$ ? The short answer is that they do and this is shown later but for now here is a wild assumption which glosses over these issues.

**Wild Assumption 7.3.1** *For every  $b > 0$  there exists a unique differentiable function  $\exp_b(x) \equiv b^x$  valid for all real values of  $x$  such that (7.4) and (7.5) both hold for all  $x, y \in \mathbb{R}$ ,  $\exp_b(m/n) = \sqrt[n]{b^m}$  whenever  $m, n$  are integers, and  $b^x > 0$  for all  $x \in \mathbb{R}$ . Furthermore, if  $b \neq 1$  and  $h \neq 0$ , then  $\exp_b(h) = b^h \neq 1$ .*

Instead of writing  $\exp_b(x)$  I will often write  $b^x$  and I will also be somewhat sloppy and regard  $b^x$  as the name of a function and not just as  $\exp_b(x)$ , a given function defined at  $x$ . This is done to conform with usual usage. Also, the last claim in Wild Assumption 7.3.1 follows from the first part of this assumption. See Problem 1. I want it to be completely clear that the Wild Assumption is just that. No reason for believing in such an assumption has been given notwithstanding the pretty pictures drawn by the calculator. Later in the book, the wild assumption will be completely justified. Based on Wild Assumption 7.3.1 one can easily find out all about  $b^x$ .

**Theorem 7.3.2** *Let  $\exp_b$  be defined in Wild Assumption 7.3.1 for  $b > 0$ . Then there exists a unique number, denoted by  $\ln b$  for  $b > 0$  satisfying*

$$\exp'_b(x) = \ln b \exp_b(x). \quad (7.6)$$

Furthermore,

$$\ln(ab) = \ln(a) + \ln(b), \quad \ln 1 = 0, \quad (7.7)$$

and for all  $y \in \mathbb{R}$ ,

$$\ln(b^y) = y \ln b. \quad (7.8)$$

$$\ln\left(\frac{a}{b}\right) = \ln(a) - \ln(b) \quad (7.9)$$

The function,  $x \rightarrow \ln x$  is differentiable and defined for all  $x > 0$  and

$$\ln'(x) = \frac{1}{x}. \quad (7.10)$$

The function,  $x \rightarrow \ln x$  is one to one on  $(0, \infty)$ . Also,  $\ln$  maps  $(0, \infty)$  onto  $(-\infty, \infty)$ .

**Proof:** First consider (7.6).

$$\begin{aligned} \lim_{h \rightarrow 0} \frac{\exp_b(x+h) - \exp_b(x)}{h} &= \lim_{h \rightarrow 0} \frac{b^{x+h} - b^x}{h} \\ &= \lim_{h \rightarrow 0} \left( \frac{b^h - 1}{h} \right) b^x. \end{aligned}$$

The expression,  $\lim_{h \rightarrow 0} \left( \frac{b^h - 1}{h} \right)$  is assumed to exist thanks to Wild Assumption 7.3.1 and this is denoted by  $\ln b$ . This proves (7.6).

To verify (7.7), if  $b = 1$  then  $b^x = 1^x$  for all  $x \in \mathbb{R}$ . Now by (7.4) and (7.5),

$$1^x 1^x = 1^{x+x} = (1^2)^x = 1^x$$

and so, dividing both sides by  $1^x$ , an operation justified by Wild Assumption 7.3.1,  $1^x = 1$  for all  $x \in \mathbb{R}$ . Therefore,  $\exp_1(x) = 1$  for all  $x$  and so  $\exp'_1(x) = \ln 1 \exp_1(x) = 0$ . Thus  $\ln 1 = 0$  as claimed. Next, by the product rule and Wild Assumption 7.3.1,

$$\begin{aligned} \ln(ab)(ab)^x &= ((ab)^x)' \\ &= (a^x b^x)' = (a^x)' b^x + a^x (b^x)' \\ &= (\ln a) a^x b^x + (\ln b) b^x a^x \\ &= [\ln a + \ln b] (ab)^x. \end{aligned}$$

Therefore,  $\ln(ab) = \ln a + \ln b$  as claimed.

Next consider (7.8). Keeping  $y$  fixed, consider the function  $x \rightarrow b^{xy} = (b^y)^x$ . Then,

$$\ln(b^y)(b^y)^x = ((b^y)^x)' = g'(x)$$

where  $g(x) \equiv b^{xy}$ . Therefore, using (6.7) on Page 122,

$$g'(x) = y \exp'_b(xy)$$

and so

$$\ln(b^y)(b^y)^x = y \exp'_b(xy) = y(\ln b)(b^{xy}) = y(\ln b)(b^y)^x.$$

Now dividing both sides by  $(b^y)^x$  verifies (7.8).

To obtain (7.9) from this, note

$$\ln\left(\frac{a}{b}\right) = \ln(ab^{-1}) = \ln(a) + \ln(b^{-1}) = \ln(a) - \ln(b).$$

It remains to verify (7.10). From (6.2) and the continuity of  $2^x$ ,

$$\ln'(1) = \lim_{h \rightarrow 0} \frac{\ln(2^h) - \ln 1}{2^h - 1} = \lim_{h \rightarrow 0} \frac{h \ln 2}{2^h - 1} = \frac{\ln 2}{\ln 2} = 1.$$



$\ln 2 \neq 0$  because if it were, then  $(2^x)' = (\ln 2) 2^x = 0$  and by Corollary 6.8.4, this would imply  $2^x$  is a constant function which it is not. Now the first part of this lemma implies

$$\begin{aligned}\ln'(x) &= \lim_{y \rightarrow x} \frac{\ln y - \ln x}{y - x} = \lim_{y \rightarrow x} \frac{1}{x} \frac{\ln\left(\frac{y}{x}\right) - \ln 1}{\left(\frac{y}{x}\right) - 1} \\ &= \frac{1}{x} \ln'(1) = \frac{1}{x}.\end{aligned}$$

It remains to verify  $\ln$  is one to one. Suppose  $\ln x = \ln y$ . Then by the mean value theorem, there exists  $t$  between  $x$  and  $y$  such that  $(1/t)(x - y) = \ln x - \ln y = 0$ . Therefore,  $x = y$  and this shows  $\ln$  is one to one as claimed.

It only remains to verify that  $\ln$  maps  $(0, \infty)$  onto  $(-\infty, \infty)$ . By Wild Assumption 7.3.1 and the mean value theorem, Corollary 6.8.3, there exists  $y \in (0, 1)$  such that

$$0 < \frac{2^1 - 1}{1} = \ln(2) 2^y$$

Since  $2^y > 0$  it follows  $\ln(2) > 0$  and so  $x \rightarrow 2^x$  is strictly increasing. Therefore, by Corollary 6.8.3

$$\frac{2^1 - 1}{1} = \ln(2) 2^y \leq \ln(2) 2^1$$

and it follows that

$$\frac{1}{2} \leq \ln 2. \quad (7.11)$$

Also, from (7.7)

$$0 = \ln(2) + \ln\left(\frac{1}{2}\right) \geq \frac{1}{2} + \ln\left(\frac{1}{2}\right)$$

which shows that

$$\ln\left(\frac{1}{2}\right) \leq -\frac{1}{2}. \quad (7.12)$$

It follows from (7.11) and (7.12) that  $\ln$  achieves values which are arbitrarily large and arbitrarily large in the negative direction. Therefore, by the intermediate value theorem,  $\ln$  achieves all values.

More precisely, let  $y \in \mathbb{R}$ . Then choose  $n$  large enough that  $\frac{n}{2} > y$  and  $-\frac{n}{2} < y$ . Then from (7.11) and (7.12)

$$\ln\left(\left(\frac{1}{2}\right)^n\right) \leq \frac{-n}{2} < y < \frac{n}{2} < \ln(2^n).$$

By the intermediate value theorem, there exists  $x \in \left(\left(\frac{1}{2}\right)^n, 2^n\right)$  such that  $\ln x = y$ . This proves the theorem.

### 7.3.3 The Special Number, $e$

Since  $\ln$  is one to one onto  $\mathbb{R}$ , it follows there exists a unique number,  $e$  such that  $\ln(e) = 1$ . Therefore,

$$\exp_e'(x) \equiv (e^x)' = \ln(e) e^x \equiv \ln(e) \exp_e(x) = \exp_e(x)$$

showing that  $\exp_e$  has the remarkable property that it equals its own derivative. This wonderful number is called Euler's number and it can be shown to equal approximately 2.7183.

### 7.3.4 The Function $\ln |x|$

The function,  $\ln$  is only defined on positive numbers. However, it is possible to write  $\ln |x|$  whenever  $x \neq 0$ . What is the derivative of this function?

**Corollary 7.3.3** *Let  $f(x) = \ln |x|$  for  $x \neq 0$ . Then*

$$f'(x) = \frac{1}{x}.$$

**Proof:** If  $x > 0$  the formula is just (7.10). Suppose then that  $x < 0$ . Then  $\ln |x| = \ln(-x)$  so by (6.7) on Page 122,

$$\begin{aligned} (\ln |x|)' &= (\ln(-x))' = (\ln((-1)x))' \\ &= \frac{1}{-x}(-1) = \frac{1}{x}. \end{aligned}$$

This proves the corollary.

### 7.3.5 Logarithm Functions

Next a new function called  $\log_b$  will be defined.

**Definition 7.3.4** *For all  $b > 0$  and  $b \neq 1$*

$$\log_b(x) \equiv \frac{\ln x}{\ln b}. \quad (7.13)$$

Notice this definition implies (7.7) - (7.9) all hold with  $\ln$  replaced with  $\log_b$ .

The fundamental relationship between the exponential function,  $b^x$  and  $\log_b x$  is in the following proposition. This proposition shows this new function is  $\log_b$  you may have studied in high school.

**Proposition 7.3.5** *Let  $b > 0$  and  $b \neq 1$ . Then for all  $x > 0$ ,*

$$b^{\log_b x} = x, \quad (7.14)$$

*and for all  $y \in \mathbb{R}$ ,*

$$\log_b b^y = y, \quad (7.15)$$

*Also,*

$$\log'_b(x) = \frac{1}{\ln b} \frac{1}{x}. \quad (7.16)$$

**Proof:** Formula (7.14) follows from (7.8).

$$\ln(b^{\log_b x}) = \log_b x \ln b = \ln x$$

and so, since  $\ln$  is one to one, it follows (7.14) holds.

$$\log_b b^y \equiv \frac{\ln(b^y)}{\ln b} = \frac{y \ln b}{\ln b} = y$$

and this verifies (7.15). Formula (7.16) is obvious from (7.13).

The functions,  $\log_b$  are only defined on positive numbers. However, it is possible to write  $\log_b |x|$  whenever  $x \neq 0$ . What is the derivative of these functions?

**Corollary 7.3.6** Let  $f(x) = \log_b |x|$  for  $x \neq 0$ . Then

$$f'(x) = \frac{1}{(\ln b)x}.$$

**Proof:** If  $x > 0$  the formula is just (7.16). Suppose then that  $x < 0$ . Then  $\log_b |x| = \log_b (-x)$  so by (6.7) on Page 122,

$$\begin{aligned} (\log_b |x|)' &= (\log_b (-x))' = (\log_b ((-1)x))' \\ &= \frac{1}{-x \ln b} (-1) = \frac{1}{x \ln b}. \end{aligned}$$

This proves the corollary.

**Example 7.3.7** Using properties of logarithms, simplify the expression,  $\log_3 \left(\frac{1}{9}x\right)$ .

From (7.7) - (7.9),

$$\begin{aligned} \log_3 \left(\frac{1}{9}x\right) &= \log_3 \left(\frac{1}{9}\right) + \log_3 (x) \\ &= \log_3 (3^{-2}) + \log_3 (x) = -2 + \log_3 (x). \end{aligned}$$

**Example 7.3.8** Using properties of logarithms, solve  $5^{x-1} = 3^{2x+2}$ .

Take  $\ln$  of both sides. Thus  $(x-1)\ln 5 = (2x+2)\ln 3$ . Then solving this for  $x$  yields  $x = \frac{\ln 5 + 2 \ln 3}{\ln 5 - 2 \ln 3}$ .

**Example 7.3.9** Solve  $\log_3 (x) + 2 = \log_9 (x+3)$ .

From the given equation,

$$3^{\log_3 (x)+2} = 3^{\log_9 (x+3)} = 9^{\frac{1}{2}(\log_9 (x+3))} = 9^{\log_9 \sqrt{x+3}}$$

and so  $9x = \sqrt{x+3}$ . Therefore,  $x = \frac{1+\sqrt{1+12 \times 81}}{2(81)} = \frac{1}{162} + \frac{1}{162}\sqrt{973}$ . In the use of the quadratic formula, only one solution was possible. (Why?)

**Example 7.3.10** Compare  $\ln(x)$  and  $\log_e(x)$

Recall that

$$\log_e(x) \equiv \frac{\ln x}{\ln e}.$$

Since  $\ln e = 1$  from the definition of  $e$ , it follows  $\log_e(x) = \ln x$ . These logarithms are called natural logarithms.

## 7.4 Exercises

1. Prove the last part of the Wild Assumption follows from the first part of this assumption. That is, show that if  $b \neq 1$ , then  $\exp_b(h) \neq 1$  if  $h \neq 0$  follows from the first part. **Hint:** If  $b^h = 1$  for  $h \neq 0$ , show  $b^x = 1$  for all  $x \in \mathbb{R}$ .

2. Simplify

(a)  $\log_4(16x)$ .

(b)  $\log_3(27x^3)$

- (c)  $(\log_b a)(\log_a b)$  for  $a, b$  positive real numbers not equal to 1.
3. Explain why the function  $6^x(1 - x \ln 6)$  is never larger than 1. **Hint:** Consider  $f(x) = 6^x(1 - x \ln 6)$  and find its maximum value.
  4. Solve  $\log_2(x) + 3 = \log_2(3x + 8)$ .
  5. Solve  $\log_4(x) + 3 = \log_2(x + 8)$ .
  6. Solve the equation  $5^{2x+9} = 7^x$  in terms of logarithms.
  7. Using properties of logarithms, simplify the expression,  $\log_4\left(\frac{1}{64}x\right)$ .
  8. Using properties of logarithms, solve  $4^{x-1} = 3^{2x+2}$ .
  9. The Wild Assumption gave the existence of a function,  $b^x$  satisfying certain properties. Show there can be no more than one such function. **Hint:** Recall the rational numbers were dense in  $\mathbb{R}$  and so one can obtain a rational number arbitrarily close to a given real number. Exploit this and the assumed continuity of  $\exp_b$  to obtain uniqueness.
  10. Prove the function,  $b^x$  is concave up and  $\log_b(x)$  is concave down.
  11. Prove that  $\log_b : (0, \infty) \rightarrow \mathbb{R}$  is onto. **Hint:** You know it is differentiable so it is continuous (Why?). Now show from (7.8) that it assumes values which are large and negative and values which are large and positive. Then use the intermediate value theorem to fill in the gaps.
  12. Using properties of logarithms and exponentials, solve  $3 + \ln(-3x) = 4 + \ln 3x^2$ .
  13. Let  $f$  be a differentiable function and suppose  $f(a) \geq 0$  and that  $f'(x) \geq 0$  for  $x \geq a$ . Show that  $f(x) \geq f(a)$  for all  $x \geq a$ . **Hint:** Use the mean value theorem.
  14. Let  $e$  be defined in Problem ?? and suppose  $e < x < y$ . Find a relationship between  $x^y$  and  $y^x$ . **Hint:** Use Problem 13 and at some point consider the function  $h(x) = \frac{\ln x}{x}$ .
  15. Suppose  $f$  is any function defined on the positive real numbers and  $f'(x) = g(x)$  where  $g$  is an odd function. ( $g(-x) = -g(x)$ .) Show  $(f(|x|))' = g(x)$ .

# Properties And Applications Of Derivatives

## 8.1 The Chain Rule And Derivatives Of Inverse Functions

### 8.1.1 The Chain Rule

The chain rule is one of the most important of differentiation rules. Special cases of it are in Theorem 6.2.6. Now it is time to consider the theorem in full generality.

**Theorem 8.1.1** Suppose  $f : (a, b) \rightarrow (c, d)$  and  $g : (c, d) \rightarrow \mathbb{R}$ . Also suppose that  $f'(x)$  exists and that  $g'(f(x))$  exists. Then  $(g \circ f)'(x)$  exists and

$$(g \circ f)'(x) = g'(f(x)) f'(x).$$

**Proof:** Define

$$H(h) \equiv \begin{cases} \frac{g(f(x+h)) - g(f(x))}{f(x+h) - f(x)} & \text{if } f(x+h) - f(x) \neq 0 \\ g'(f(x)) & \text{if } f(x+h) - f(x) = 0 \end{cases}.$$

Then for  $h \neq 0$ ,

$$\frac{g(f(x+h)) - g(f(x))}{h} = H(h) \frac{f(x+h) - f(x)}{h}.$$

Note that  $\lim_{h \rightarrow 0} H(h) = g'(f(x))$  due to Theorems 5.9.6 on Page 100 and 6.2.2 on Page 121. Therefore, taking the limit and using Theorem 5.9.4,

$$\lim_{h \rightarrow 0} \frac{g(f(x+h)) - g(f(x))}{h} = g'(f(x)) f'(x).$$

This proves the chain rule.

**Example 8.1.2** Let  $f(x) = \ln |\ln(x^4 + 1)|$ . Find  $f'(x)$ .

From the chain rule,

$$\begin{aligned} f'(x) &= \ln'(\ln(x^4 + 1)) (\ln(x^4 + 1))' \\ &= \frac{1}{\ln(x^4 + 1)} \frac{1}{x^4 + 1} (x^4)' \\ &= \left( \frac{1}{\ln(x^4 + 1)} \right) \left( \frac{1}{x^4 + 1} \right) (4x^3). \end{aligned}$$

**Example 8.1.3** Let  $f(x) = (2 + \ln|x|)^3$ . Find  $f'(x)$ .

Use the chain rule again. Thus

$$\begin{aligned} f'(x) &= 3(2 + \ln|x|)^2 (2 + \ln|x|)' \\ &= \frac{3}{x} (2 + \ln|x|)^2. \end{aligned}$$

### 8.1.2 Implicit Differentiation And Derivatives Of Inverse Functions

Sometimes a function is not given explicitly in terms of a formula. For example, you might have  $x^2 + y^2 = 4$ . This relation defines  $y$  as a function of  $x$  near a given point such as  $(0, 1)$ . Near this point,  $y = \sqrt{4 - x^2}$ . Near the point,  $(0, -1)$ , you have  $y = -\sqrt{4 - x^2}$ . Near the point,  $(1, 0)$ , you can't solve for  $y$  in terms of  $x$  but you can solve for  $x$  in terms of  $y$ . Thus near  $(1, 0)$ ,  $x = \sqrt{4 - y^2}$ . This was a simple example but in general, you can't use algebra to solve for one of the variables in terms of the others even if the relation defines that variable as a function of the others. Here is an example in which, even though it is impossible to find  $y(x)$  you can still find the derivative of  $y$ . The procedure by which this is accomplished is nothing more than the chain rule and other rules of differentiation.

**Example 8.1.4** Suppose  $y$  is a differentiable function of  $x$  and  $y^3 + 2yx = x^3 + 7 + \ln|y|$ . Find  $y'(x)$ .

This illustrates the technique of implicit differentiation. If you believe  $y$  is some differentiable function of  $x$ , then you can differentiate both sides with respect to  $x$  and write, using the chain rule and product rule.

$$3y^2y' + 2xy' + 2y = 3x^2 + \frac{y'}{y}.$$

Now you can solve for  $y'$  and obtain  $y' = -\frac{2y-3x^2}{3y^3+2xy-1}y$ .

Of course there are significant mathematical considerations which are being ignored when it is assumed  $y$  is a differentiable function of  $x$ . It turns out that for problems like this, the equation relating  $x$  and  $y$  actually does define  $y$  as a differentiable function of  $x$  near points where it makes sense to formally solve for  $y'$  as just done. The theorems which give this justification are called the implicit and inverse function theorems. They are some of the most profound theorems in mathematics and are topics for advanced calculus. The interested reader should consult the book by Rudin, [14] for this and generalizations of all the hard theorems given in this book. One case is of special interest in which  $y = f(x)$  and it is desired to find  $\frac{dx}{dy}$  or in other words, the derivative of the inverse function.

It happens that if  $f$  is a differentiable one to one function defined on an interval,  $[a, b]$ , and  $f'(x)$  exists and is non zero then the inverse function,  $f^{-1}$  has a derivative at the point  $f(x)$ . Recall that  $f^{-1}$  is defined according to the formula

$$f^{-1}(f(x)) = x.$$

**Definition 8.1.5** Let  $f : [a, b] \rightarrow \mathbb{R}$  be a continuous function. Define

$$f'(a) \equiv \lim_{x \rightarrow a+} \frac{f(x) - f(a)}{x - a}, \quad f'(b) \equiv \lim_{x \rightarrow b-} \frac{f(x) - f(b)}{x - b}.$$

Recall the notation  $x \rightarrow a+$  means that only  $x > a$  are considered in the definition of limit. The notation  $x \rightarrow b-$  is defined similarly. Thus, this definition includes the derivative of  $f$  at the endpoints of the interval and to save notation,

$$f'(x_1) \equiv \lim_{x \rightarrow x_1} \frac{f(x) - f(x_1)}{x - x_1}$$

where it is understood that  $x$  is always in  $[a, b]$ .

**Theorem 8.1.6** *Let  $f : [a, b] \rightarrow \mathbb{R}$  be continuous and one to one. Suppose  $f'(x_1)$  exists for some  $x_1 \in [a, b]$  and  $f'(x_1) \neq 0$ . Then  $(f^{-1})'(f(x_1))$  exists and is given by the formula,  $(f^{-1})'(f(x_1)) = \frac{1}{f'(x_1)}$ .*

**Proof:** By Lemma 5.7.4, and Corollary 5.7.6 on Page 95  $f$  is either strictly increasing or strictly decreasing and  $f^{-1}$  is continuous. Therefore there exists  $\eta > 0$  such that if  $0 < |f(x_1) - f(x)| < \eta$ , then

$$0 < |x_1 - x| = |f^{-1}(f(x_1)) - f^{-1}(f(x))| < \delta$$

where  $\delta$  is small enough that for  $0 < |x_1 - x| < \delta$ ,

$$\left| \frac{x - x_1}{f(x) - f(x_1)} - \frac{1}{f'(x_1)} \right| < \varepsilon.$$

It follows that if  $0 < |f(x_1) - f(x)| < \eta$ ,

$$\left| \frac{f^{-1}(f(x)) - f^{-1}(f(x_1))}{f(x) - f(x_1)} - \frac{1}{f'(x_1)} \right| = \left| \frac{x - x_1}{f(x) - f(x_1)} - \frac{1}{f'(x_1)} \right| < \varepsilon$$

Therefore, since  $\varepsilon > 0$  is arbitrary,

$$\lim_{y \rightarrow f(x_1)} \frac{f^{-1}(y) - f^{-1}(f(x_1))}{y - f(x_1)} = \frac{1}{f'(x_1)}$$

and this proves the theorem.

The following obvious corollary comes from the above by not bothering with end points.

**Corollary 8.1.7** *Let  $f : (a, b) \rightarrow \mathbb{R}$  be continuous and one to one. Suppose  $f'(x_1)$  exists for some  $x_1 \in (a, b)$  and  $f'(x_1) \neq 0$ . Then  $(f^{-1})'(f(x_1))$  exists and is given by the formula,  $(f^{-1})'(f(x_1)) = \frac{1}{f'(x_1)}$ .*

This is one of those theorems which is very easy to remember if you neglect the difficult questions and simply focus on formal manipulations. Consider the following.

$$f^{-1}(f(x)) = x.$$

Now use the chain rule on both sides to write

$$(f^{-1})'(f(x)) f'(x) = 1,$$

and then divide both sides by  $f'(x)$  to obtain

$$(f^{-1})'(f(x)) = \frac{1}{f'(x)}.$$

Of course this gives the conclusion of the above theorem rather effortlessly and it is formal manipulations like this which aid many of us in remembering formulas such as the one given in the theorem.

**Example 8.1.8** Let  $f(x) = \ln(1+x^2) + x^3 + 7$ . Show that  $f$  has an inverse and find  $(f^{-1})'(7)$ .

I am not able to find a formula for the inverse function. This is typical in useful applications so you need to get used to this idea. The methods of algebra are insufficient to solve hard problems in analysis. You need something more. The question is to determine whether  $f$  has an inverse. To do this,

$$\begin{aligned} f'(x) &= \frac{2x}{1+x^2} + 3x^2 + 7 \\ &> -1 + 3x^2 + 7 \\ &> 6 > 0. \end{aligned}$$

By Corollary 6.8.5 on Page 133, this function is strictly increasing on  $\mathbb{R}$  and so it has an inverse function although I have no idea how to find an explicit formula for this inverse function. However, I can see that  $f(0) = 7$  and so by the formula for the derivative of an inverse function,

$$\begin{aligned} (f^{-1})'(7) &= (f^{-1})'(f(0)) = \frac{1}{f'(0)} \\ &= \frac{1}{7}. \end{aligned}$$

**Example 8.1.9** Suppose  $f(a) = 0$  and  $f'(x) = \sqrt{1+x^4 + \ln(1+x^2)}$ . Find  $(f^{-1})'(0)$ .

The function,  $f$  is one to one because it is strictly increasing due to the fact that its derivative is positive for all  $x$ . As in the last example, I have no idea how to find a formula for  $f^{-1}$  but I do see that  $f(a) = 0$  and so

$$(f^{-1})'(0) = (f^{-1})'(f(a)) = \frac{1}{f'(a)} = \frac{1}{\sqrt{1+a^4 + \ln(1+a^2)}}.$$

The chain rule has a particularly attractive form in Leibniz's notation. Suppose  $y = g(u)$  and  $u = f(x)$ . Thus  $y = g \circ f(x)$ . Then from the above theorem

$$\begin{aligned} (g \circ f)'(x) &= g'(f(x)) f'(x) \\ &= g'(u) f'(x) \end{aligned}$$

or in other words,

$$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx}.$$

Notice how the  $du$ 's cancel. This particular form is a very useful crutch and is used extensively in applications.

## 8.2 Exercises

1. In each of the following, find  $\frac{dy}{dx}$ .

(a)  $y = e^{\sin x}$

(b)  $y = \ln(\sin(x^2 + 7))$

(c)  $y = \tan(\cos(x^2))$



- (d)  $y = \log_2 (\sin(x) + 6)$
  - (e)  $y = \sin (\log_3 (x^2 + 1))$
  - (f)  $y = \frac{\sqrt{x^3+7}}{\sqrt{\sin(x)+4}}$
  - (g)  $y = 3^{\tan(\sin(x))}$
  - (h)  $y = \left( \frac{x^2+2x}{\tan(x^2+1)} \right)^6$
2. In each of the following, assume the relation defines  $y$  as a function of  $x$  for values of  $x$  and  $y$  of interest and use the process of implicit differentiation to find  $y'(x)$ .
- (a)  $xy^2 + \sin(y) = x^3 + 1$
  - (b)  $y^3 + x \cos(y^2) = x^4$
  - (c)  $y \cos(x) = \tan(y) \cos(x^2) + 2$
  - (d)  $(x^2 + y^2)^6 = x^3y + 3$
  - (e)  $\frac{xy^2+y}{y^5+x} + \cos(y) = 7$
  - (f)  $\sqrt{x^2 + y^4} \sin(y) = 3x$
  - (g)  $y^3 \sin(x) + y^2x^2 = 2^{x^2}y + \ln|y|$
  - (h)  $y^2 \sin(y)x + \log_3(xy) = y^2 + 11$
  - (i)  $\sin(x^2 + y^2) + \sec(xy) = e^{x+y} + y2^y + 2$
  - (j)  $\sin(\tan(xy^2)) + y^3 = 16$
  - (k)  $\cos(\sec(\tan(y))) + \ln(5 + \sin(xy)) = x^2y + 3$
3. Show that if  $D(g) \subseteq U \subseteq D(f)$ , and if  $f$  and  $g$  are both one to one, then  $f \circ g$  is also one to one.
4. Using Problem 3 show that the following functions are one to one and find the derivative of the inverse function at the indicated point.
- (a)  $y = e^{x^3+1}, e^2$
  - (b)  $y = (x^3 + 7x + 1)^3, 0$
  - (c)  $y = \tan(x^3 + \frac{\pi}{4}), 1$
  - (d)  $y = \tan(-x^5 + \frac{\pi}{4}), 1$
  - (e)  $y = 2^{5x+\sin(x)}, 5(\frac{\pi}{2}) + 1$

## 8.3 The Function $x^r$ For $r$ A Real Number

Theorem 7.3.2 on Page 143 says that for  $x > 0$ , and for  $r$  a real number,

$$\ln(x^r) = r \ln(x) \quad (8.1)$$

By this theorem, it also follows that  $\ln^{-1} : \mathbb{R} \rightarrow (0, \infty)$  exists. Then by Corollary 8.1.7  $\ln^{-1}$  is differentiable. This function is so important it is given a special symbol,  $\exp$ . Thus

$$\exp(\ln x) = x, \ln(\exp(y)) = y.$$

**Proposition 8.3.1** *There exists a unique number  $e$  such that  $\ln e = 1$ . Then  $\exp(x) = e^x$ . Furthermore,  $\exp'(x) = \exp(x)$ .*

**Proof:** From Theorem 7.3.2 on Page 143 there exists a unique number,  $e$  such that  $\ln e = 1$ . Now by this theorem again,  $\ln(e^x) = x \ln e = x$ . Also, from the above,  $\ln(\exp x) = x$  and so since  $\ln$  is one to one, it follows  $e^x = \exp x$  as claimed.

To establish the last claim, note  $\ln(\exp x) = x$  and so using the chain rule and Corollary 8.1.7, it follows

$$\frac{\exp'(x)}{\exp(x)} = 1$$

which gives the desired result.

Note that

$$x = \ln(\exp(x)), \quad x = x \ln(e) = \ln(e^x) = \ln(\exp_e(x))$$

and so, since  $\ln$  is one to one, it follows  $\exp(x) = \exp_e(x) = e^x$ .

With this understanding, it becomes possible to find derivatives of functions raised to arbitrary real powers. First, note that (8.1) can be written as

$$x^r = \exp(r \ln x) \tag{8.2}$$

**Theorem 8.3.2** *For  $x > 0$ ,  $(x^r)' = rx^{r-1}$ .*

**Proof:** Differentiate both sides of (8.2) using the chain rule. From the Wild Assumption on Page 143, in particular, the part about the validity of the laws of exponents,

$$(x^r)' = \exp'(r \ln x) \frac{r}{x} = \exp(r \ln x) \frac{r}{x} = \frac{r}{x} x^r = rx^{r-1}. \tag{8.3}$$

and this shows from (8.3) that  $(x^r)' = rx^{r-1}$  as claimed.

**Example 8.3.3** *Suppose  $f(x)$  is a non zero differentiable function. Find the derivative of  $|f(x)|^r$ .*

From (8.2),

$$|f(x)|^r = \exp(r \ln |f(x)|).$$

Therefore,

$$\begin{aligned} (|f(x)|^r)' &= \exp(r \ln |f(x)|) (r \ln |f(x)|)' \\ &= |f(x)|^r r \frac{f'(x)}{f(x)} = r |f(x)|^{r-2} f(x) f'(x). \end{aligned}$$

### 8.3.1 Logarithmic Differentiation

**Example 8.3.4** *Let  $f(x) = (1 + x^2)^x$ . Find  $f'(x)$ .*

One way to do this is to take  $\ln$  of both sides and use the chain rule to differentiate both sides with respect to  $x$ . Thus

$$\ln(f(x)) = x \ln(1 + x^2)$$

and so, taking the derivative of both sides, using the chain and product rules,

$$\frac{f'(x)}{f(x)} = \frac{2x^2}{1 + x^2} + \ln(1 + x^2).$$

Then solve for  $f'(x)$  to obtain

$$f'(x) = (1+x^2)^x \left( \frac{2x^2}{1+x^2} + \ln(1+x^2) \right).$$

This process is called logarithmic differentiation.

**Example 8.3.5** Let  $f(x) = \frac{\sqrt[3]{x^3 + \sin(x)}}{\sqrt[6]{x^4 + 2x}}$ . Find  $f'(x)$ .

You could use the quotient and chain rules but it is easier to use logarithmic differentiation.

$$\ln(f(x)) = \frac{1}{3} \ln(x^3 + \sin(x)) - \frac{1}{6} \ln(x^4 + 2x).$$

Differentiating both sides,

$$\frac{f'(x)}{f(x)} = \frac{1}{3} \left( \frac{3x^2 + \cos x}{x^3 + \sin x} \right) - \frac{1}{6} \frac{4x^3 + 2}{x^4 + 2x}.$$

Therefore, the answer is

$$f'(x) = \frac{\sqrt[3]{x^3 + \sin(x)}}{\sqrt[6]{x^4 + 2x}} \left( \frac{1}{3} \left( \frac{3x^2 + \cos x}{x^3 + \sin x} \right) - \frac{1}{6} \frac{4x^3 + 2}{x^4 + 2x} \right).$$

I think you can see the advantage of doing it this way over using the quotient rule.

## 8.4 Exercises

1. Let  $f(x) \equiv x^3 + 1$ . Find  $f^{-1}(y)$ . Now find  $(f^{-1})'(1)$ .
2. Let  $f(x) \equiv x^3 + 7x + 3$ . Explain why  $f$  has an inverse. Find  $(f^{-1})'(3)$ .
3. Derive the quotient rule from the product rule and the chain rule. This shows you don't need to remember the wretched quotient rule if you don't want to. It follows from two rules which you cannot survive without.
4. What is wrong with the following "proof" of the chain rule? Here  $g'(f(x))$  exists and  $f'(x)$  exists.

$$\begin{aligned} & \lim_{h \rightarrow 0} \frac{g(f(x+h)) - g(f(x))}{h} \\ &= \lim_{h \rightarrow 0} \frac{g(f(x+h)) - g(f(x))}{f(x+h) - f(x)} \frac{f(x+h) - f(x)}{h} \\ &= g'(f(x)) f'(x). \end{aligned}$$

5. Is the derivative of a function always continuous? **Hint:** Consider a differentiable function which is periodic of period 1 and non constant,  $f$ . Periodic of period 1 means  $f(x+1) = f(x)$  for all  $x \in \mathbb{R}$ . Now consider

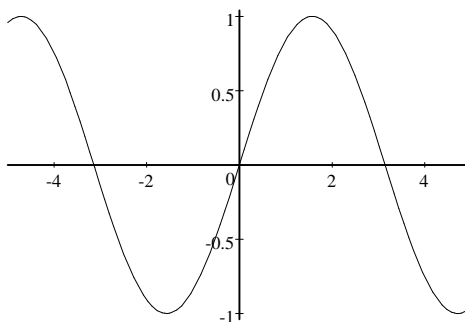
$$h(x) = \begin{cases} x^2 f\left(\frac{1}{x}\right) & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases}.$$

Show  $h'(0) = 0$ . What is  $h'(x)$  for  $x \neq 0$ ? Is  $h'$  also periodic of period 1?

6. Let  $f(x) = x^3 + 1$ . Find an explicit formula for  $f^{-1}$  and use it to compute  $(f^{-1})'(9)$ . Then use the formula given in the theorem of this section to see you get the same answer.
7. Find the derivatives of the following functions.
- (a)  $\sin(x^2) \ln(x^2 + 1)$
  - (b)  $\ln(1 + x^2)$
  - (c)  $(x^3 + 1)^6 \sin(x^2 + 7)$
  - (d)  $\ln((x^3 + 1)^6 \sin(x^2 + 7))$
  - (e)  $\tan(\sec(\sin(x^2 + 1)))$
  - (f)  $(\sin^2(x^2 + 5))^{\sqrt{7}}$
8. Use (8.2) or logarithmic differentiation to differentiate the following functions.
- (a)  $(2 + \sin(x^2 + 6))^{\tan x}$
  - (b)  $x^x$
  - (c)  $(x^x)^x$
  - (d)  $(\tan^2(x^4 + 4) + 1)^{\cos x}$
  - (e)  $(\sin^2(x))^{\tan x}$

## 8.5 The Inverse Trigonometric Functions

It is desired to consider the inverse trigonometric functions. Graphing the function  $y = \sin x$ , is clear  $\sin$  is not one to one on  $\mathbb{R}$  and so it is not possible to define an inverse function.



However, a little thing like this will not prevent the definition of useful inverse trig. functions. Observe the function,  $\sin$ , is one to one on the interval  $[-\frac{\pi}{2}, \frac{\pi}{2}]$  shown in the above picture as the interval containing zero on which the function climbs from  $-1$  to  $1$ . Also note that for  $y$  on this interval,  $\cos(y) \geq 0$ . Now the arcsin function is defined as the inverse of the  $\sin$  when its domain is restricted to  $[-\frac{\pi}{2}, \frac{\pi}{2}]$ . In words,  $\arcsin(x)$  is defined to be the angle whose sine is  $x$  which lies in  $[-\frac{\pi}{2}, \frac{\pi}{2}]$ . From Theorem 8.1.6 on Page 151 about the derivative of the inverse function, the derivative of  $x \rightarrow \arcsin(x)$  exists for all  $x \in [-1, 1]$ . The formula in this theorem could be used to find the derivative of  $\arcsin$  but it is more useful to simply use that theorem to resolve the existence question and apply

the chain rule to find the formula. It is a mistake to memorize too many formulas. Let  $y = \arcsin(x)$  so  $\sin(y) = x$ . Now taking the derivative of both sides,

$$\cos(y) y' = 1$$

and so

$$y' = \frac{1}{\cos y} = \frac{1}{\sqrt{1 - \sin^2(y)}} = \frac{1}{\sqrt{1 - x^2}}.$$

The positive value of the square root is used because for  $y \in [-\frac{\pi}{2}, \frac{\pi}{2}]$ ,  $\cos(y) \geq 0$ . Thus

$$\frac{1}{\sqrt{1 - x^2}} = \arcsin'(x). \quad (8.4)$$

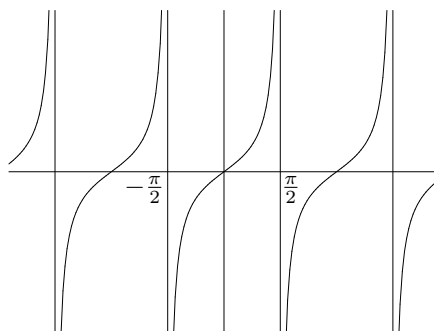
Next consider the inverse tangent function. You are aware that  $\tan$  is periodic of period  $\pi$  because

$$\begin{aligned} \tan(x + \pi) &= \frac{\sin(x + \pi)}{\cos(x + \pi)} = \frac{\sin(x)\cos(\pi) + \cos(x)\sin(\pi)}{\cos(x)\cos(\pi) - \sin(x)\sin(\pi)} \\ &= \frac{-\sin(x)}{-\cos(x)} = \tan(x). \end{aligned}$$

Therefore, it is impossible to take the inverse of  $\tan$ . However,  $\tan$  is one to one on  $(-\frac{\pi}{2}, \frac{\pi}{2})$  and

$$\lim_{x \rightarrow \frac{\pi}{2}^-} \tan(x) = +\infty, \quad \lim_{x \rightarrow -\frac{\pi}{2}^+} \tan(x) = -\infty$$

as shown in the following graph of  $y = \tan(x)$  in which the vertical lines represent vertical asymptotes.



Therefore,  $\arctan(x)$  for  $x \in (-\infty, \infty)$  is defined according to the rule:  $\arctan(x)$  is the angle whose tangent is  $x$  which is in  $(-\frac{\pi}{2}, \frac{\pi}{2})$ . By Theorem 8.1.6  $\arctan$  has a derivative. Therefore, letting  $y = \arctan(x)$ ,  $\tan(y) = x$  and by the chain rule,

$$\sec^2(y) y' = 1.$$

Therefore,

$$y' = \frac{1}{\sec^2(y)} = \frac{1}{1 + \tan^2(y)} = \frac{1}{1 + x^2}.$$

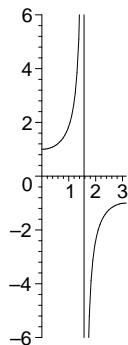
and so

$$\frac{1}{1 + x^2} = \arctan'(x). \quad (8.5)$$

The inverse secant function can be defined similarly. There is no agreement on the best way to restrict the domain of  $\sec$ . I will follow the way of doing it which is used in the book

by Salas and Hille [16] recognizing that there are good reasons for doing it other ways also. The graph of  $\sec$  is represented below on  $[0, \frac{\pi}{2}) \cup (\frac{\pi}{2}, \pi]$ . There is a vertical asymptote at  $x = \frac{\pi}{2}$ . Thus

$$\lim_{x \rightarrow \frac{\pi}{2}^-} \sec(x) = +\infty, \quad \lim_{x \rightarrow \frac{\pi}{2}^+} \sec(x) = -\infty$$



As in the case of  $\arcsin$  and  $\arctan$ ,  $\operatorname{arcsec}(x)$  is the angle whose secant is  $x$  which lies in  $[0, \frac{\pi}{2}) \cup (\frac{\pi}{2}, \pi]$ . Let  $y = \operatorname{arcsec}(x)$  so  $x = \sec(y)$  and using the chain rule,

$$1 = \sec(y) \tan(y) y'.$$

Now from the trig. identity,  $1 + \tan^2(y) = \sec^2(y)$ ,

$$\begin{aligned} y' &= \frac{1}{\sec(y) \tan(y)} \\ &= \frac{1}{x (\pm \sqrt{x^2 - 1})} \end{aligned}$$

and it is necessary to consider what to do with  $\pm$ . If  $y \in [0, \frac{\pi}{2})$ , both  $x = \sec(y)$  and  $\tan(y)$  are nonnegative and so in this case,

$$y' = \frac{1}{x\sqrt{x^2 - 1}} = \frac{1}{|x|\sqrt{x^2 - 1}}.$$

If  $y \in (\frac{\pi}{2}, \pi]$ , then  $x = \sec(y) < 0$  and  $\tan(y) \leq 0$  so

$$y' = \frac{1}{x(-1)\sqrt{x^2 - 1}} = \frac{1}{(-x)\sqrt{x^2 - 1}} = \frac{1}{|x|\sqrt{x^2 - 1}}.$$

Thus either way,

$$y' = \frac{1}{|x|\sqrt{x^2 - 1}}.$$

This yields the formula

$$\frac{1}{|x|\sqrt{x^2 - 1}} = \operatorname{arcsec}'(x). \quad (8.6)$$

As in the case of  $\ln$ , there is an interesting and useful formula involving  $\operatorname{arcsec}(|x|)$ . For  $x < 0$ , this function equals  $\operatorname{arcsec}(-x)$  and so by the chain rule, its derivative equals

$$\frac{1}{|-x|\sqrt{x^2 - 1}}(-1) = \frac{1}{x\sqrt{x^2 - 1}}$$

If  $x > 0$ , this function equals  $\operatorname{arcsec}(x)$  and its derivative equals

$$\frac{1}{|x|\sqrt{x^2-1}} = \frac{1}{x\sqrt{x^2-1}}$$

and so either way,

$$(\operatorname{arcsec}(|x|))' = \frac{1}{x\sqrt{x^2-1}}. \quad (8.7)$$

## 8.6 The Hyperbolic And Inverse Hyperbolic Functions

The hyperbolic functions are given by

$$\sinh(x) \equiv \frac{e^x - e^{-x}}{2}, \cosh(x) \equiv \frac{e^x + e^{-x}}{2},$$

and

$$\tanh(x) \equiv \frac{\sinh(x)}{\cosh(x)}.$$

The first of these is called the hyperbolic sine and the second the hyperbolic cosine. I imagine you can guess what the third is called. If you guessed “hyperbolic tangent” you got it right. The other hyperbolic functions are defined by analogy to the circular functions.

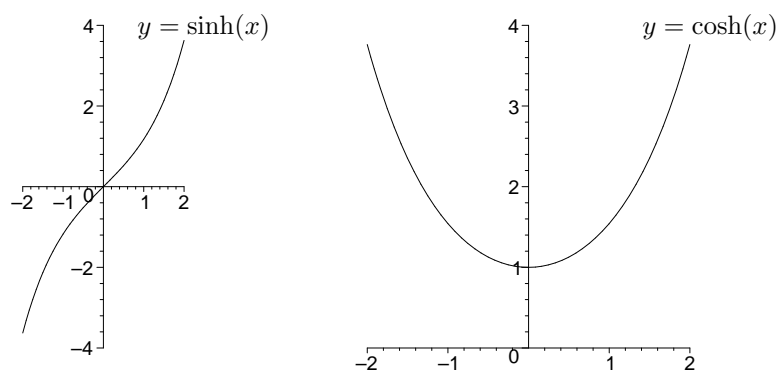
The reason these are called hyperbolic functions is that

$$\cosh^2 t - \sinh^2 t = 1 \quad (8.8)$$

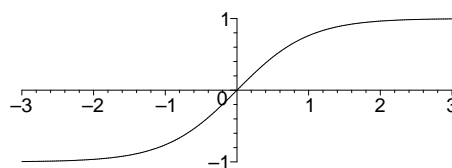
and so the point,  $(\cosh t, \sinh t)$  is a point on the hyperbola whose equation is  $x^2 - y^2 = 1$ . This is not important but is the source for the term hyperbolic. Using the chain rule,

$$\cosh'(x) = \sinh(x), \sinh'(x) = \cosh x.$$

Also, you see that  $\sinh(0) = 0$ ,  $\cosh(0) = 1$  and that  $\sinh(x) < 0$  if  $x < 0$  while  $\sinh(x) > 0$  for  $x > 0$ , but  $\cosh(x) > 0$  for all  $x$ . Therefore,  $\sinh$  is an increasing function, concave down for  $x < 0$  and concave up for  $x > 0$  because  $\sinh''(x) = \sinh(x)$  while  $\cosh$  is decreasing for  $x < 0$  and increasing for  $x > 0$ . Since  $\cosh''(x) = \cosh(x)$ , it is concave up for all  $x$ . Thus the graphs of these functions are as follows.



Also, you can use the graph of the function  $x \rightarrow \exp(x) = e^x$  to verify that the graph of  $\tanh(x)$  is as given below



Since  $x \rightarrow \sinh(x)$  is strictly increasing, it has an inverse function,  $\sinh^{-1}(x)$ . If  $y = \sinh^{-1}(x)$ , then  $\sinh(y) = x$  and so using the chain rule and the theorem about the existence of the derivative of the inverse function,  $y' \cosh(y) = 1$ . From the identity (8.8)  $\cosh(y) = \sqrt{1 + \sinh^2(y)} = \sqrt{1 + x^2}$ . Therefore,

$$y' = \frac{1}{\sqrt{1 + x^2}}$$

which gives the formula

$$\frac{1}{\sqrt{1 + x^2}} = (\sinh^{-1})'(x). \quad (8.9)$$

The derivative of the hyperbolic tangent is also easy to find. This yields after a short computation

$$(\tanh x)' = 1 - \tanh^2 x. \quad (8.10)$$

Another notation for the inverse hyperbolic functions which is sometimes used is  $\operatorname{arcsinh}$  or  $\operatorname{arccosh}$  or  $\operatorname{arctanh}$ .

## 8.7 Exercises

1. Verify (8.10).
2. Find derivatives of the following functions.
  - (a)  $\sinh^{-1}(x^2 + 7)$
  - (b)  $\tanh^{-1}(x)$
  - (c)  $\sin(\sinh^{-1}(x^2 + 2))$
  - (d)  $\sin(\tanh(x))$
  - (e)  $x^2 \sinh(\sin(\cos(x)))$
  - (f)  $(\cosh(x^3))^{\sqrt{6}}$
  - (g)  $(1 + x^4)^{\sin x}$
3. Simplify  $\arcsin x + \arccos x$ .
4. A wonderful identity which was used to compute  $\pi$  for over 200 years<sup>1</sup> is the following.

$$\frac{\pi}{4} = 4 \arctan\left(\frac{1}{5}\right) - \arctan\left(\frac{1}{239}\right).$$

<sup>1</sup>John Machin computed  $\pi$  to 100 decimal places in 1706 through the use of this identity. Later in 1873 William Shanks did it to over 700 places using this identity. The next advance was in 1948, 808 decimal places. After this, computers began to be used and currently  $\pi$  is “known” to millions of decimal places. Many other schemes have been used besides this identity for computing  $\pi$ .



Establish this identity by taking the tangent of both sides and using an appropriate formula for the tangent of the difference of two angles. Use De Moivre's theorem to get some help in finding a formula for  $\tan(4\theta)$ .

5. Find a formula for  $\tanh^{-1}$  in terms of  $\ln$ .
6. Find a formula for  $\sinh^{-1}$  in terms of  $\ln$ .
7. Prove  $1 - \tanh^2 x = \operatorname{sech}^2 x$ .
8. Prove  $\coth^2 x - 1 = \operatorname{csch}^2 x$ .
9. What about  $\cosh^{-1}$ ? Define it by restricting the domain of  $\cosh$  to be nonnegative numbers? What is  $\cosh^{-1}$  in terms of  $\ln$ ?
10. Show  $\arcsin x = \arctan\left(\frac{x}{\sqrt{1-x^2}}\right)$ . It is possible to start with the  $\arctan$  function and obtain all the other trig functions in terms of this one. If you knew the function,  $\arctan$  explain how to define  $\sin$  and  $\cos$ . This is interesting because there is a simple way to define  $\arctan$  directly as a function of a real variable[9]. Approaches like these avoid all reference to plane geometry.

## 8.8 L'Hôpital's Rule

There is an interesting rule which is often useful for evaluating difficult limits called L'Hôpital's<sup>2</sup> rule. The best versions of this rule are based on the Cauchy Mean value theorem, Theorem 6.8.2 on Page 133.

**Theorem 8.8.1** *Let  $[a, b] \subseteq [-\infty, \infty]$  and suppose  $f, g$  are functions which satisfy,*

$$\lim_{x \rightarrow b-} f(x) = \lim_{x \rightarrow b-} g(x) = 0, \quad (8.11)$$

*and  $f'$  and  $g'$  exist on  $(a, b)$  with  $g'(x) \neq 0$  on  $(a, b)$ . Suppose also that*

$$\lim_{x \rightarrow b-} \frac{f'(x)}{g'(x)} = L. \quad (8.12)$$

*Then*

$$\lim_{x \rightarrow b-} \frac{f(x)}{g(x)} = L. \quad (8.13)$$

**Proof:** By the definition of limit and (8.12) there exists  $c < b$  such that if  $t > c$ , then

$$\left| \frac{f'(t)}{g'(t)} - L \right| < \frac{\varepsilon}{2}.$$

Now pick  $x, y$  such that  $c < x < y < b$ . By the Cauchy mean value theorem, there exists  $t \in (x, y)$  such that

$$g'(t)(f(x) - f(y)) = f'(t)(g(x) - g(y)).$$

---

<sup>2</sup>L'Hôpital published the first calculus book in 1696. This rule, named after him, appeared in this book. The rule was actually due to Bernoulli who had been L'Hôpital's teacher. L'Hôpital did not claim the rule as his own but Bernoulli accused him of plagiarism. Nevertheless, this rule has become known as L'Hôpital's rule ever since. The version of the rule presented here is superior to what was discovered by Bernoulli and depends on the Cauchy mean value theorem which was found over 100 years after the time of L'Hôpital.

Since  $g'(s) \neq 0$  for all  $s \in (a, b)$  it follows  $g(x) - g(y) \neq 0$ . Therefore,

$$\frac{f'(t)}{g'(t)} = \frac{f(x) - f(y)}{g(x) - g(y)}$$

and so, since  $t > c$ ,

$$\left| \frac{f(x) - f(y)}{g(x) - g(y)} - L \right| < \frac{\varepsilon}{2}.$$

Now letting  $y \rightarrow b-$ ,

$$\left| \frac{f(x)}{g(x)} - L \right| \leq \frac{\varepsilon}{2} < \varepsilon.$$

Since  $\varepsilon > 0$  is arbitrary, this shows (8.13).

The following corollary is proved in the same way.

**Corollary 8.8.2** *Let  $[a, b] \subseteq [-\infty, \infty]$  and suppose  $f, g$  are functions which satisfy,*

$$\lim_{x \rightarrow a+} f(x) = \lim_{x \rightarrow a+} g(x) = 0, \quad (8.14)$$

*and  $f'$  and  $g'$  exist on  $(a, b)$  with  $g'(x) \neq 0$  on  $(a, b)$ . Suppose also that*

$$\lim_{x \rightarrow a+} \frac{f'(x)}{g'(x)} = L. \quad (8.15)$$

*Then*

$$\lim_{x \rightarrow a+} \frac{f(x)}{g(x)} = L. \quad (8.16)$$

Here is a simple example which illustrates the use of this rule.

**Example 8.8.3** Find  $\lim_{x \rightarrow 0} \frac{5x + \sin 3x}{\tan 7x}$ .

The conditions of L'Hôpital's rule are satisfied because the numerator and denominator both converge to 0 and the derivative of the denominator is nonzero for  $x$  close to 0. Therefore, if the limit of the quotient of the derivatives exists, it will equal the limit of the original function. Thus,

$$\lim_{x \rightarrow 0} \frac{5x + \sin 3x}{\tan 7x} = \lim_{x \rightarrow 0} \frac{5 + 3 \cos 3x}{7 \sec^2(7x)} = \frac{8}{7}.$$

Sometimes you have to use L'Hôpital's rule more than once.

**Example 8.8.4** Find  $\lim_{x \rightarrow 0} \frac{\sin x - x}{x^3}$ .

Note that  $\lim_{x \rightarrow 0} (\sin x - x) = 0$  and  $\lim_{x \rightarrow 0} x^3 = 0$ . Also, the derivative of the denominator is nonzero for  $x$  close to 0. Therefore, if  $\lim_{x \rightarrow 0} \frac{\cos x - 1}{3x^2}$  exists and equals  $L$ , it will follow from L'Hôpital's rule that the original limit exists and equals  $L$ . However,  $\lim_{x \rightarrow 0} (\cos x - 1) = 0$  and  $\lim_{x \rightarrow 0} 3x^2 = 0$  so L'Hôpital's rule can be applied again to consider  $\lim_{x \rightarrow 0} \frac{-\sin x}{6x}$ . From L'Hôpital's rule, if this limit exists and equals  $L$ , it will follow that  $\lim_{x \rightarrow 0} \frac{\cos x - 1}{3x^2} = L$  and consequently  $\lim_{x \rightarrow 0} \frac{\sin x - x}{x^3} = L$ . But from Lemma 7.1.1 on Page 139,  $\lim_{x \rightarrow 0} \frac{-\sin x}{6x} = \frac{-1}{6}$ . Therefore, by L'Hôpital's rule,  $\lim_{x \rightarrow 0} \frac{\sin x - x}{x^3} = \frac{-1}{6}$ .

**Warning 8.8.5** *Be sure to check the assumptions of L'Hôpital's rule before using it.*

**Example 8.8.6** Find  $\lim_{x \rightarrow 0+} \frac{\cos 2x}{x}$ .

The numerator becomes close to 1 and the denominator gets close to 0. Therefore, the assumptions of L'Hôpital's rule do not hold and so it does not apply. In fact there is no limit unless you define the limit to equal  $+\infty$ . Now let's try to use the conclusion of L'Hôpital's rule even though the conditions for using this rule are not verified. Take the derivative of the numerator and the denominator which yields  $\frac{-2 \sin 2x}{1}$ , an expression whose limit as  $x \rightarrow 0+$  equals 0. This is a good illustration of the above warning.

Some people get the unfortunate idea that one can find limits by doing experiments with a calculator. If the limit is taken as  $x$  gets close to 0, these people think one can find the limit by evaluating the function at values of  $x$  which are closer and closer to 0. Theoretically, this should work although you have no way of knowing how small you need to take  $x$  to get a good estimate of the limit. In practice, the procedure may fail miserably.

**Example 8.8.7** Find  $\lim_{x \rightarrow 0} \frac{\ln|1+x^{10}|}{x^{10}}$ .

This limit equals  $\lim_{y \rightarrow 0} \frac{\ln|1+y|}{y} = \lim_{y \rightarrow 0} \frac{\left(\frac{1}{1+y}\right)}{1} = 1$  where L'Hôpital's rule has been used. This is an amusing example. You should plug .001 in to the function,  $\frac{\ln|1+x^{10}|}{x^{10}}$  and see what your calculator or computer gives you. If it is like mine, it will give the answer, 0 and will keep on returning the answer of 0 for smaller numbers than .001. This illustrates the folly of trying to compute limits through calculator or computer experiments.

There is another form of L'Hôpital's rule in which  $\lim_{x \rightarrow b-} f(x) = \pm\infty$  and  $\lim_{x \rightarrow b-} g(x) = \pm\infty$ .

**Theorem 8.8.8** Let  $[a, b] \subseteq [-\infty, \infty]$  and suppose  $f, g$  are functions which satisfy,

$$\lim_{x \rightarrow b-} f(x) = \pm\infty \text{ and } \lim_{x \rightarrow b-} g(x) = \pm\infty, \quad (8.17)$$

and  $f'$  and  $g'$  exist on  $(a, b)$  with  $g'(x) \neq 0$  on  $(a, b)$ . Suppose also

$$\lim_{x \rightarrow b-} \frac{f'(x)}{g'(x)} = L. \quad (8.18)$$

Then

$$\lim_{x \rightarrow b-} \frac{f(x)}{g(x)} = L. \quad (8.19)$$

**Proof:** By the definition of limit and (8.18) there exists  $c < b$  such that if  $t > c$ , then

$$\left| \frac{f'(t)}{g'(t)} - L \right| < \frac{\varepsilon}{2}.$$

Now pick  $x, y$  such that  $c < x < y < b$ . By the Cauchy mean value theorem, there exists  $t \in (x, y)$  such that

$$g'(t)(f(x) - f(y)) = f'(t)(g(x) - g(y)).$$

Since  $g'(s) \neq 0$  on  $(a, b)$ , it follows  $g(x) - g(y) \neq 0$ . Therefore,

$$\frac{f'(t)}{g'(t)} = \frac{f(x) - f(y)}{g(x) - g(y)}$$

and so, since  $t > c$ ,

$$\left| \frac{f(x) - f(y)}{g(x) - g(y)} - L \right| < \frac{\varepsilon}{2}.$$

Now this implies

$$\left| \frac{f(y) \left( \frac{f(x)}{f(y)} - 1 \right)}{g(y) \left( \frac{g(x)}{g(y)} - 1 \right)} - L \right| < \frac{\varepsilon}{2}$$

where for all  $y$  large enough, both  $\frac{f(x)}{f(y)} - 1$  and  $\frac{g(x)}{g(y)} - 1$  are not equal to zero. Continuing to rewrite the above inequality yields

$$\left| \frac{f(y)}{g(y)} - L \frac{\left( \frac{g(x)}{g(y)} - 1 \right)}{\left( \frac{f(x)}{f(y)} - 1 \right)} \right| < \frac{\varepsilon}{2} \left| \frac{\left( \frac{g(x)}{g(y)} - 1 \right)}{\left( \frac{f(x)}{f(y)} - 1 \right)} \right|.$$

Therefore, for  $y$  large enough,

$$\left| \frac{f(y)}{g(y)} - L \right| \leq \left| L - L \frac{\left( \frac{g(x)}{g(y)} - 1 \right)}{\left( \frac{f(x)}{f(y)} - 1 \right)} \right| + \frac{\varepsilon}{2} \left| \frac{\left( \frac{g(x)}{g(y)} - 1 \right)}{\left( \frac{f(x)}{f(y)} - 1 \right)} \right| < \varepsilon$$

due to the assumption (8.17) which implies

$$\lim_{y \rightarrow b-} \frac{\left( \frac{g(x)}{g(y)} - 1 \right)}{\left( \frac{f(x)}{f(y)} - 1 \right)} = 1.$$

Therefore, whenever  $y$  is large enough,

$$\left| \frac{f(y)}{g(y)} - L \right| < \varepsilon$$

and this is what is meant by (8.19). This proves the theorem.

As before, there is no essential difference between the proof in the case where  $x \rightarrow b-$  and the proof when  $x \rightarrow a+$ . This observation is stated as the next corollary.

**Corollary 8.8.9** *Let  $[a, b] \subseteq [-\infty, \infty]$  and suppose  $f, g$  are functions which satisfy,*

$$\lim_{x \rightarrow a+} f(x) = \pm\infty \text{ and } \lim_{x \rightarrow a+} g(x) = \pm\infty, \quad (8.20)$$

*and  $f'$  and  $g'$  exist on  $(a, b)$  with  $g'(x) \neq 0$  on  $(a, b)$ . Suppose also that*

$$\lim_{x \rightarrow a+} \frac{f'(x)}{g'(x)} = L. \quad (8.21)$$

*Then*

$$\lim_{x \rightarrow a+} \frac{f(x)}{g(x)} = L. \quad (8.22)$$

Theorems 8.8.1 8.8.8 and Corollaries 8.8.2 and 8.8.9 will be referred to as L'Hôpital's rule from now on. Theorem 8.8.1 and Corollary 8.8.2 involve the notion of indeterminate forms of the form  $\frac{0}{0}$ . Please do not think any meaning is being assigned to the nonsense expression  $\frac{0}{0}$ . It is just a symbol to help remember the sort of thing described by Theorem 8.8.1 and Corollary 8.8.2. Theorem 8.8.8 and Corollary 8.8.9 deal with indeterminate forms which are of the form  $\frac{\pm\infty}{\infty}$ . Again, this is just a symbol which is helpful in remembering the sort of thing being considered. There are other indeterminate forms which can be reduced to these forms just discussed. Don't ever try to assign meaning to such symbols.

**Example 8.8.10** Find  $\lim_{y \rightarrow \infty} \left(1 + \frac{x}{y}\right)^y$ .

It is good to first see why this is called an indeterminate form. One might think that as  $y \rightarrow \infty$ , it follows  $x/y \rightarrow 0$  and so  $1 + \frac{x}{y} \rightarrow 1$ . Now 1 raised to anything is 1 and so it would seem this limit should equal 1. On the other hand, if  $x > 0$ ,  $1 + \frac{x}{y} > 1$  and a number raised to higher and higher powers should approach  $\infty$ . It really isn't clear what this limit should be. It is an indeterminate form which can be described as  $1^\infty$ . By definition,

$$\left(1 + \frac{x}{y}\right)^y = \exp\left(y \ln\left(1 + \frac{x}{y}\right)\right).$$

Now using L'Hôpital's rule,

$$\begin{aligned} \lim_{y \rightarrow \infty} y \ln\left(1 + \frac{x}{y}\right) &= \lim_{y \rightarrow \infty} \frac{\ln\left(1 + \frac{x}{y}\right)}{1/y} \\ &= \lim_{y \rightarrow \infty} \frac{\frac{1}{1+(x/y)}(-x/y^2)}{(-1/y^2)} \\ &= \lim_{y \rightarrow \infty} \frac{x}{1+(x/y)} = x \end{aligned}$$

Therefore,

$$\lim_{y \rightarrow \infty} y \ln\left(1 + \frac{x}{y}\right) = x$$

Since  $\exp$  is continuous, it follows

$$\lim_{y \rightarrow \infty} \left(1 + \frac{x}{y}\right)^y = \lim_{y \rightarrow \infty} \exp\left(y \ln\left(1 + \frac{x}{y}\right)\right) = e^x.$$

### 8.8.1 Interest Compounded Continuously

Suppose you put money in the bank and it accrues interest at the rate of  $r$  per payment period. These terms need a little explanation. If the payment period is one month, and you started with \$100 then the amount at the end of one month would equal  $100(1+r) = 100 + 100r$ . In this the second term is the interest and the first is called the principal. Now you have  $100(1+r)$  in the bank. This becomes the new principal. How much will you have at the end of the second month? By analogy to what was just done it would equal

$$100(1+r) + 100(1+r)r = 100(1+r)^2.$$

In general, the amount you would have at the end of  $n$  months is  $100(1+r)^n$ .

When a bank says they offer 6% compounded monthly, this means  $r$ , the rate per payment period equals .06/12. Consider the problem of a rate of  $r$  per year and compounding the interest  $n$  times a year and letting  $n$  increase without bound. This is what is meant by compounding continuously. The interest rate per payment period is then  $r/n$  and the number of payment periods after time  $t$  years is approximately  $tn$ . From the above the amount in the account after  $t$  years is

$$P \left(1 + \frac{r}{n}\right)^{nt} \tag{8.23}$$

Recall from Example 8.8.10 that  $\lim_{y \rightarrow \infty} \left(1 + \frac{x}{y}\right)^y = e^x$ . The expression in (8.23) can be written as

$$P \left[\left(1 + \frac{r}{n}\right)^n\right]^t$$

and so, taking the limit as  $n \rightarrow \infty$ , you get

$$Pe^{rt} = A.$$

This shows how to compound interest continuously.

**Example 8.8.11** Suppose you have \$100 and you put it in a savings account which pays 6% compounded continuously. How much will you have at the end of 4 years?

From the above discussion, this would be  $100e^{(.06)4} = 127.12$ . Thus, in 4 years, you would gain interest of about \$27.

## 8.9 Exercises

1. Find the limits.

- (a)  $\lim_{x \rightarrow 0} \frac{3x-4 \sin 3x}{\tan 3x}$
- (b)  $\lim_{x \rightarrow \frac{\pi}{2}^-} (\tan x)^{x-(\pi/2)}$
- (c)  $\lim_{x \rightarrow 1} \frac{\arctan(4x-4)}{\arcsin(4x-4)}$
- (d)  $\lim_{x \rightarrow 0} \frac{\arctan 3x-3x}{x^3}$
- (e)  $\lim_{x \rightarrow 0^+} \frac{9^{\sec x-1}-1}{3^{\sec x-1}-1}$
- (f)  $\lim_{x \rightarrow 0} \frac{3x+\sin 4x}{\tan 2x}$
- (g)  $\lim_{x \rightarrow \pi/2} \frac{\ln(\sin x)}{x-(\pi/2)}$
- (h)  $\lim_{x \rightarrow 0} \frac{\cosh 2x-1}{x^2}$
- (i)  $\lim_{x \rightarrow 0} \frac{-\arctan x+x}{x^3}$
- (j)  $\lim_{x \rightarrow 0} \frac{x^8 \sin \frac{1}{x}}{\sin 3x}$
- (k)  $\lim_{x \rightarrow \infty} (1+5^x)^{\frac{2}{x}}$
- (l)  $\lim_{x \rightarrow 0} \frac{-2x+3 \sin x}{x}$
- (m)  $\lim_{x \rightarrow 1} \frac{\ln(\cos(x-1))}{(x-1)^2}$
- (n)  $\lim_{x \rightarrow 0^+} \sin^{\frac{1}{x}} x$
- (o)  $\lim_{x \rightarrow 0} (\csc 5x - \cot 5x)$
- (p)  $\lim_{x \rightarrow 0^+} \frac{3^{\sin x}-1}{2^{\sin x}-1}$
- (q)  $\lim_{x \rightarrow 0^+} (4x)^{x^2}$
- (r)  $\lim_{x \rightarrow \infty} \frac{x^{10}}{(1.01)^x}$
- (s)  $\lim_{x \rightarrow 0} (\cos 4x)^{(1/x^2)}$

2. Find the following limits.

- (a)  $\lim_{x \rightarrow 0^+} \frac{1-\sqrt{\cos 2x}}{\sin^4(4\sqrt{x})}$ .
- (b)  $\lim_{x \rightarrow 0} \frac{2^{x^2}-2^{5x}}{\sin\left(\frac{x^2}{5}\right)-\sin(3x)}$ .

- (c)  $\lim_{n \rightarrow \infty} n \left( \sqrt[n]{7} - 1 \right).$
- (d)  $\lim_{x \rightarrow \infty} \left( \frac{3x+2}{5x-9} \right)^{x^2}.$
- (e)  $\lim_{x \rightarrow \infty} \left( \frac{3x+2}{5x-9} \right)^{1/x}.$
- (f)  $\lim_{n \rightarrow \infty} \left( \cos \frac{2x}{\sqrt{n}} \right)^n.$
- (g)  $\lim_{n \rightarrow \infty} \left( \cos \frac{2x}{\sqrt{5n}} \right)^n.$
- (h)  $\lim_{x \rightarrow 3} \frac{x^x - 27}{x - 3}.$
- (i)  $\lim_{n \rightarrow \infty} \cos \left( \pi \frac{\sqrt{4n^2 + 13n}}{n} \right).$
- (j)  $\lim_{x \rightarrow \infty} \left( \sqrt[3]{x^3 + 7x^2} - \sqrt{x^2 - 11x} \right).$
- (k)  $\lim_{x \rightarrow \infty} \left( \sqrt[5]{x^5 + 7x^4} - \sqrt[3]{x^3 - 11x^2} \right).$
- (l)  $\lim_{x \rightarrow \infty} \left( \frac{5x^2 + 7}{2x^2 - 11} \right)^{\frac{x}{1-x}}.$
- (m)  $\lim_{x \rightarrow \infty} \left( \frac{5x^2 + 7}{2x^2 - 11} \right)^{\frac{x \ln x}{1-x}}.$
- (n)  $\lim_{x \rightarrow 0+} \frac{\ln(e^{2x^2} + 7\sqrt{x})}{\sinh(\sqrt{x})}.$
- (o)  $\lim_{x \rightarrow 0+} \frac{\sqrt[7]{x} - \sqrt[5]{x}}{\sqrt[9]{x} - \sqrt[11]{x}}.$

3. Find the following limits.

- (a)  $\lim_{x \rightarrow 0+} (1 + 3x)^{\cot 2x}$
- (b)  $\lim_{x \rightarrow 0} \frac{\tan(\sin x) - \sin(\tan x)}{x^7}$
- (c)  $\lim_{x \rightarrow 0} \frac{\sin(x^2) - \sin^2(x)}{x^4}$
- (d)  $\lim_{x \rightarrow 0} \frac{e^{-(1/x^2)}}{x}$
- (e)  $\lim_{x \rightarrow 0} \left( \frac{1}{x} - \cot(x) \right)$
- (f)  $\lim_{x \rightarrow 0} \frac{\cos(\sin x) - 1}{x^2}$
- (g)  $\lim_{x \rightarrow \infty} \left( x^2 (4x^4 + 7)^{1/2} - 2x^4 \right)$
- (h)  $\lim_{x \rightarrow 0} \frac{\cos(x) - \cos(4x)}{\tan(x^2)}$
- (i)  $\lim_{x \rightarrow 0} \frac{\arctan(3x)}{x}$
- (j)  $\lim_{x \rightarrow \infty} \left[ (x^9 + 5x^6)^{1/3} - x^3 \right]$

## 8.10 Related Rates

Sometimes some variables are related by a formula and it is known how fast all are changing but one. The related rates problem asks for how fast the remaining variable is changing.

**Example 8.10.1** A cube of ice is melting such that  $\frac{dV}{dt} = -4\text{cm}^3/\text{sec}$  where  $V$  is the volume. How fast are the sides changing when they are equal to 5 centimeters in length?

The volume is  $V = x^3$  where  $x$  is the length of a side of the cube. Therefore, the chain rule implies

$$-4 = \frac{dV}{dt} = 3x^2 \frac{dx}{dt}$$

and the problem is to find  $\frac{dx}{dt}$  when  $x = 5$ . Therefore,

$$\frac{dx}{dt} = \frac{-4}{3(25)} = \frac{-4}{75} \text{ cm/second}$$

at this time.

Note there is no way of knowing the volume or the sides as a functions of  $t$ .

**Example 8.10.2** *One car travels north at 70 miles per hour and the other travels east at 60 miles per hour toward an intersection. How fast is the distance between the two cars changing when this distance equals five miles and the car heading north is at a distance of three miles from the intersection?*

Let  $l$  denote the distance between the cars. Thus if  $x$  is the distance from the intersection of the car traveling east and  $y$  is the distance from the intersection of the car traveling north,  $l^2 = x^2 + y^2$ . When  $y = 3$ , it follows that  $x = 4$ . Therefore, at the instant described

$$\begin{aligned} 2ll' &= 2xx' + 2yy' \\ 10l' &= 8(-60) + 6(-70) \end{aligned}$$

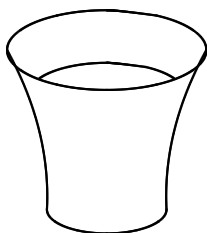
and so  $l' = -90$  miles per hour at this instant.

## 8.11 Exercises

1. One car travels north at 70 miles per hour toward an intersection and the other travels east at 60 miles per hour away from the intersection. How fast is the distance between the two cars changing when this distance equals five miles and the car heading north is at a distance of three miles from the intersection?
2. A trash compactor compacts some trash which is in the shape of a box having a square base and a height equal to twice the length of a side of the base. Suppose each side of the base is changing at the rate of  $-3$  inches per second. How fast is the volume changing when the side of the base equals 10 inches?
3. An isosceles triangle has two sides of equal length. Imagine such a triangle in which the two legs have length 8 inches and denote the included angle by  $\theta$  and the area by  $A$ . Suppose  $\frac{dA}{dt} = \sqrt{3}$  square inches per minute. How fast is  $\theta$  changing when  $\theta = \pi/6$  radians?
4. A point having coordinates  $(x, y)$  moves over the ellipse,  $\frac{x^2}{4} + \frac{y^2}{9} = 1$ . If  $\frac{dy}{dt} = 2$ , find  $\frac{dx}{dt}$  at the point  $(2, 3)$ .
5. A spectator at a tennis tournament sits 10 feet from the end of the net and on the line determined by the net. He watches the ball go back and forth, and will have a very sore neck when he wakes up the next morning. How fast is the angle between his line of sight and the line determined by the net changing when the ball crosses over the net at a point 12 feet from the end assuming the ball travels at a speed of 60 miles per hour?



6. The surface area of a sphere of radius  $r$  equals  $4\pi r^2$  and the volume of the ball of radius  $r$  equals  $(4/3)\pi r^3$ . A balloon in the shape of a ball is being inflated at the rate of 6 cubic inches per minute. How fast is the surface area changing when the volume of the ball equals  $20\pi$  cubic inches?
7. A mother cheetah attempts to fix dinner, a Thompson gazelle, for her hungry children. She moves at 100 feet per second while dinner travels at 80 feet per second. How fast is the distance between her and dinner decreasing when she is located at the point  $(0, 40)$  feet and dinner is moving in the direction of the positive  $x$  axis at the point  $(30, 0)$  feet? Is your answer a little surprising? **Hint:** Let  $(x, y)$  denote the coordinates of the cheetah and let  $(z, 0)$  denote the coordinates of the gazelle. Then if  $l$  is the desired distance,  $l^2 = (x - z)^2 + y^2$ . At the instant described, assuming the cheetah moves toward the gazelle at all times,  $y'/x' = -4/3$  (why?) and also  $\sqrt{(x')^2 + (y')^2} = 100$ .
8. A six foot high man walks at a speed of 5 feet per second away from a light post which is 12 feet high that has the light right on the top. How fast is the end of his shadow moving when he is at a distance of 10 feet from the base of the light pole. How fast is the end of the shadow moving when he is 5 feet from the pole? (Assume he does not walk normally but instead oozes along like a giant amoeba so that his head is always exactly 6 feet above the ground.)
9. The volume of a right circular cone is  $\frac{1}{3}\pi r^2 h$ . Grain comes off a conveyor belt and falls to the ground making a right circular cone. It is observed that  $r'(t) = .5$  feet per minute and  $h'(t) = .3$  feet per minute. It is also known that the rate at which the grain falls off the conveyor belt is  $100\pi$  cubic feet per minute. When the radius of the cone is 10 feet what is the height of the cone?
10. A hemispherical dish of radius 5 inches is sitting on a table. Soup is being poured in at the constant rate of 4 cubic inches per second. How fast is the level of soup rising when the radius of the top surface of the soup equals 3 inches? The volume of soup at depth  $y$  will be shown later to equal  $V(y) = \pi \left( 5y^2 - \frac{y^3}{3} \right)$ .
11. A vase of water is sitting on a table.



It will be shown later that if  $V(y)$  is the total volume of the vase up to height  $y$ , then  $\frac{dV}{dy} = A(y)$  where  $A(y)$  is the surface area of the top surface of the water at this height. (To see this is very reasonable, note that a little chunk of volume of the vase between heights  $y$  and  $y + dy$  would be  $dV = A(y) dy$ , area times height.) Also, the rate at which the water evaporates is proportional to the surface area of the exposed water. Thus  $\frac{dV}{dt} = -kA(y)$ . Show  $\frac{dy}{dt}$  is a constant even though the surface area of the exposed water is constantly changing in a typical vase.

12. A revolving search light at a prison makes one revolution per minute. How fast is the light travelling along the nearest point on a wall  $1/4$  mile away? Give your answer in miles per hour.

13. A painter is on top of a 13 foot ladder which leans against a house. The base of the ladder is moving away from the house at the rate of 2 feet per second causing the top of the ladder to move down the house. How fast is the painter descending when the base of the ladder is at a distance of 5 feet from the house?
14. A rope fastened to the bow of a row boat has the other end wound around a windlass which is 4 feet above the level of the bow of the boat. The current pulls the boat away at the rate of 2 feet per second. How fast is the rope unwinding when the distance between the bow of the boat and the windlass is 10 feet?
15. A kite 100 feet above the ground is being blown away from the person holding its string in a direction parallel to the ground at the rate of 10 feet per second. At what rate must the string be let out when the length of string already let out is 200 feet?
16. A certain volume of an ideal gas satisfies  $PV = kT$  where  $T$  is the absolute temperature,  $P$  is the pressure,  $V$  is the volume and  $k$  is a constant which depends on the amount of the gas and the sort of gas in the sample. Find a formula for  $\frac{dV}{dt}$  in terms of  $k, P, V$  and their derivatives.
17. A disposable cup is made in the shape of a right circular cone with height 5 inches and radius 2 inches. Water flows in to this conical cup at the rate of 4 cubic inches per minute. How fast is the water level rising when the water in the cone is three inches deep? The volume of a cone is  $\frac{1}{3}\pi r^2 h$ .
18. The two equal sides of an isosceles triangle have length  $x$  inches and the third leg has length  $y$  inches. Suppose  $\frac{dx}{dt} = 2$  inches per minute and that the length of the other side changes in such a way that the area of the triangle is always 10 square inches. For  $\theta$  the angle between the two equal sides, find  $\frac{d\theta}{dt}$  when  $x = 5$  inches.

## 8.12 The Derivative And Optimization

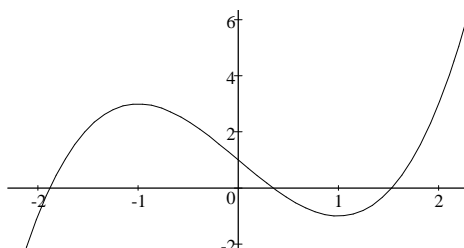
There are existence theorems such as Theorem 5.7.10 on Page 96 which ensure a maximum or minimum of a function exists but in this section the goal is to give ways to find the maximum or minimum values of a function.

Suppose  $f$  is continuous on  $[a, b]$ . The minimum or maximum could occur at either end point or it could occur at a point in the open interval,  $(a, b)$ . If it occurs at a point of the open interval,  $(a, b)$ , say at  $x_0$ , and if  $f'(x_0)$  exists, then from Theorem 6.5.2 on Page 126  $f'(x_0) = 0$ . Therefore, the following simple procedure can be used to locate the maximum or minimum of a function,  $f$ . Find all points,  $x$ , in  $(a, b)$  where  $f'(x) = 0$  and all points,  $x$ , in  $(a, b)$  where  $f'(x)$  does not exist. Then consider these points along with the end points of the interval. Evaluate  $f$  at the end points and at these points where the derivative is zero or does not exist. The largest must be the maximum value of  $f$  on the interval,  $[a, b]$ , and the smallest must be the minimum value of  $f$  on the interval,  $[a, b]$ . Typically, this involves checking only finitely many points.

Sometimes there are no end points. In this case, you do not necessarily know a maximum value or a minimum value of a continuous function even exists. However, if the function is differentiable and if a maximum or minimum exists, it can still be found by looking at the points where the derivative equals zero.

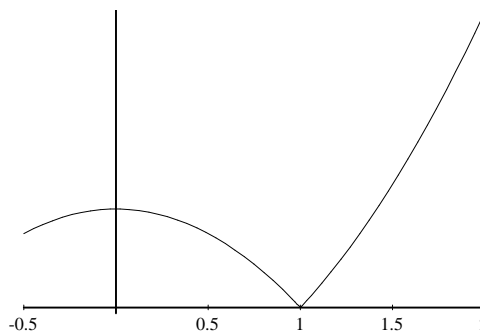
**Example 8.12.1** Find the maximum and minimum values of the function,  $f(x) = x^3 - 3x + 1$  on the interval  $[-2, 2]$ .

The points where  $f'(x) = 3x^2 - 3 = 0$  are  $x = 1$  or  $-1$ . There are no points where the derivative does not exist. Therefore, evaluate the function at  $-1, -2, 2$ , and  $1$ . Thus  $f(-1) = 3, f(1) = -1, f(-2) = -1$ , and  $f(2) = 3$ . Therefore, the maximum value of the function occurs at the point  $-1$  and  $2$  and has the value of  $3$  while the minimum value of the function occurs at  $-1$  and  $-2$  and equals  $-1$ . The following is a graph of this function.



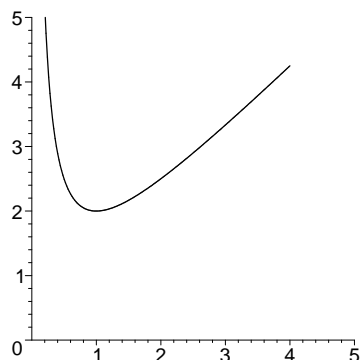
**Example 8.12.2** Find the maximum and minimum values of the function  $f(x) = |x^2 - 1|$  on the interval  $[-.5, 2]$ .

You should verify that this function fails to have a derivative at the point  $x = 1$ . For  $x \in (-.5, 1)$  the function equals  $1 - x^2$  and so its derivative equals zero at  $x = 0$ . For  $x > 1$ , the function equals  $x^2 - 1$  and so there is no point larger than  $1$  where the derivative equals zero. Therefore, the points to look at are the end points,  $-.5, 2$ , the points where the derivative fails to exist,  $1$  and the point where the derivative equals zero,  $x = 0$ . Now  $f(-.5) = .75, f(0) = 1, f(1) = 0$ , and  $f(2) = 3$ . It follows the function achieves its maximum at the end point,  $x = 2$  and its minimum at the point,  $x = 1$  where the derivative fails to exist. The following is a graph of this function.



**Example 8.12.3** Find the minimum value of the function  $f(x) = x + \frac{1}{x}$  for  $x \in (0, \infty)$ .

The graph of this function is given below.



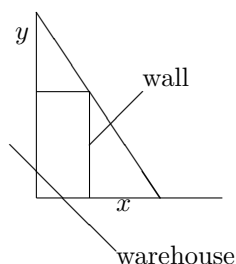
From the graph, it seems there should exist a minimum value at the bottom of the graph. To find it, take the derivative of  $f$  and set it equal to zero and then solve for the value of  $x$ . Thus

$$1 - \frac{1}{x^2} = 0$$

and so  $x = 1$ . The solution to the equation,  $x = -1$  is of no interest because it is not greater than zero. Therefore, the minimum value of the function on the interval,  $(0, \infty)$  equals  $f(1) = 2$  as suggested by the graph.

**Example 8.12.4** *An eight foot high wall stands one foot from a warehouse. What is the length of the shortest ladder which extends from the ground to the warehouse.*

A diagram of this situation is the following picture.



In this picture, the slanted line represents the ladder and  $x$  and  $y$  are as shown. By similar triangles,  $y/1 = 8/x$ . Therefore,  $xy = 8$ . From the Pythagorean theorem the length of the ladder equals  $\sqrt{1 + y^2} + \sqrt{x^2 + 64}$ . Now using the relation between  $x$  and  $y$ , the function of a single variable,  $x$ , to minimize is

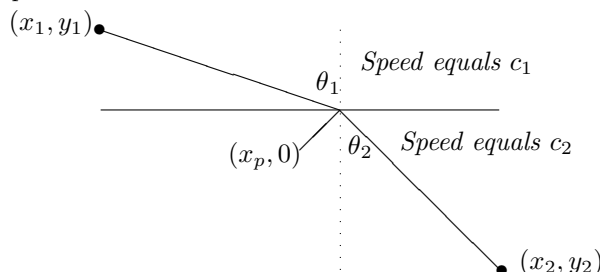
$$f(x) = \sqrt{1 + \frac{64}{x^2}} + \sqrt{x^2 + 64},$$

Clearly  $x > 0$  so there are no endpoints to worry about. (Why?) Also the function is differentiable and so it suffices to consider only points where the derivative equals zero. This is a little messy but finally

$$f'(x) = \frac{1}{\sqrt{(x^2 + 64)}} \frac{-64 + x^3}{x^2}$$

and the value where this equals zero is  $x = 4$ . It follows the shortest ladder is of length  $\sqrt{1 + \frac{64}{4^2}} + \sqrt{4^2 + 64} = 5\sqrt{5}$  feet.

**Example 8.12.5** *Fermat's principle says that light travels on a path which will minimize the total time. Consider the following picture of light passing from  $(x_1, y_1)$  to  $(x_2, y_2)$  as shown. The angle  $\theta_1$  is called the angle of incidence while the angle,  $\theta_2$  is called the angle of refraction. The picture indicates a situation in which  $c_1 > c_2$ .*



Then define by  $x$  the quantity  $x_p - x_1$ . What is the relation between  $\theta_1$  and  $\theta_2$ ?

The time it takes for the light to go from  $(x_1, y_1)$  to the point  $(x_p, 0)$  equals  $\sqrt{x^2 + y_1^2}/c_1$  and the time it takes to go from  $(x_p, 0)$  to  $(x_2, y_2)$  is  $\sqrt{(x_2 - x_1 - x)^2 + y_2^2}/c_2$ . Therefore, the total time is

$$T = \frac{\sqrt{x^2 + y_1^2}}{c_1} + \frac{\sqrt{(x_2 - x_1 - x)^2 + y_2^2}}{c_2}$$

Thus  $T$  is minimized if

$$\begin{aligned} \frac{dT}{dx} &= \frac{d}{dx} \left( \frac{\sqrt{x^2 + y_1^2}}{c_1} + \frac{\sqrt{(x_2 - x_1 - x)^2 + y_2^2}}{c_2} \right) \\ &= \frac{x}{c_1 \sqrt{x^2 + y_1^2}} - \frac{(x_2 - x_1 - x)}{c_2 \sqrt{(x_2 - x_1 - x)^2 + y_2^2}} \\ &= \frac{\sin(\theta_1)}{c_1} - \frac{\sin(\theta_2)}{c_2} = 0 \end{aligned}$$

at this point. Therefore, this yields the desired relation between  $\theta_1$  and  $\theta_2$ . This is called Snell's law.

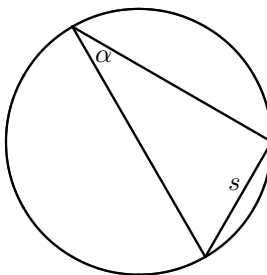
## 8.13 Exercises

1. A cylindrical can is to be constructed to hold 30 cubic inches. The top and bottom of the can are constructed of a material costing one cent per square inch and the sides are constructed of a material costing 2 cents per square inch. Find the minimum cost for such a can.
2. Two positive numbers sum to 8. Find the numbers if their product is to be as large as possible.
3. The ordered pair,  $(x, y)$  is on the ellipse,  $x^2 + 4y^2 = 4$ . Form the rectangle which has  $(x, y)$  as one end of a diagonal and  $(0, 0)$  at the other end. Find the rectangle of this sort which has the largest possible area.
4. A rectangle is inscribed in a circle of radius  $r$ . Find the formula for the rectangle of this sort which has the largest possible area.

5. A point is picked on the ellipse,  $x^2 + 4y^2 = 4$  which is in the first quadrant. Then a line tangent to this point is drawn which intersects the  $x$  axis at a point,  $x_1$  and the  $y$  axis at the point  $y_1$ . The area of the triangle formed by the  $y$  axis, the  $x$  axis, and the line just drawn is thus  $\frac{x_1 y_1}{2}$ . Out of all possible triangles formed in this way, find the one with smallest area.
6. Find maximum and minimum values if they exist for the function,  $f(x) = \frac{\ln x}{x}$  for  $x > 0$ .
7. Describe how you would find the maximum value of the function,  $f(x) = \frac{\ln x}{2 + \sin x}$  for  $x \in (0, 6)$  if it exists. **Hint:** You might want to use a calculator to graph this and get an idea what is going on.
8. A rectangular beam of height  $h$  and width  $w$  is to be sawed from a circular log of radius 1 foot. Find the dimensions of the strongest such beam assuming the strength is of the form  $kh^2w$ . Here  $k$  is some constant which depends on the type of wood used.
9. A farmer has 600 feet of fence with which to enclose a rectangular piece of land that borders a river. If he can use the river as one side, what is the largest area that he can enclose.
10. An open box is to be made by cutting out little squares at the corners of a rectangular piece of cardboard which is 20 inches wide and 40 inches long. and then folding up the rectangular tabs which result. What is the largest possible volume which can be obtained?
11. A feeding trough is to be made from a rectangular piece of metal which is 3 feet wide and 12 feet long by folding up two rectangular pieces of dimension one foot by 12 feet. What is the best angle for this fold?
12. Find the dimensions of the right circular cone which has the smallest area given the volume is  $30\pi$  cubic inches. The volume of the right circular cone is  $(1/3)\pi r^2 h$  and the area of the cone is  $\pi r \sqrt{h^2 + r^2}$ .
13. A wire of length 10 inches is cut into two pieces, one of length  $x$  and the other of length  $10 - x$ . One piece is bent into the shape of a square and the other piece is bent into the shape of a circle. Find the two lengths such that the sum of the areas of the circle and the square is as large as possible. What are the lengths if the sum of the two areas is to be as small as possible.
14. A hiker begins to walk to a cabin in a dense forest. He is walking on a road which runs from East to West and the cabin is located exactly one mile north of a point two miles down the road. He walks 5 miles per hour on the road but only 3 miles per hour in the woods. Find the path which will minimize the time it takes for him to get to the cabin.
15. A refinery is on a straight shore line. Oil needs to flow from a mooring place for oil tankers to this refinery. Suppose the mooring place is two miles off shore from a point on the shore 8 miles away from the refinery and that it costs five times as much to lay pipe under water than above the ground. Describe the most economical route for a pipeline from the mooring place to the refinery.
16. Two hallways, one 5 feet wide and the other 6 feet wide meet. It is desired to carry a ladder horizontally around the corner. What is the longest ladder which can be carried in this way? **Hint:** Consider a line through the inside corner which extends

to the opposite walls. The shortest such line will be the length of the longest ladder. You might also consider Example 7.1.3 on Page 141.

17. A triangle is inscribed in a circle in such a way that one side of the triangle is always the same length,  $s$ . Show that out of all such triangles the maximum area is obtained when the triangle is an isosceles triangle. **Hint:** From theorems in plane geometry, the angle opposite the side having fixed length is a constant, no matter how you draw the triangle. Use the law of sines and this fact. In the following picture,  $\alpha$  is a constant.

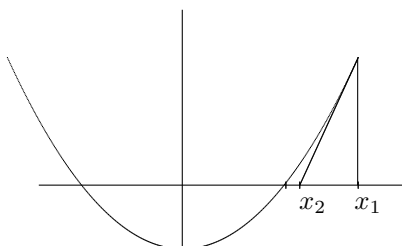


18. A window is to be constructed for the wall of a church which is to consist of a rectangle of height  $b$  surmounted by a half circle of radius  $a$ . Suppose the total perimeter of the window is to be no more than  $4\pi + 8$  feet. Find the shape and dimensions of the window which will admit the most light.
19. You know  $\lim_{x \rightarrow \infty} \ln x = \infty$ . Show that if  $\alpha > 0$ , then  $\lim_{x \rightarrow \infty} \frac{\ln x}{x^\alpha} = 0$ .
20. Suppose  $p$  and  $q$  are two positive numbers larger than 1 which satisfy  $\frac{1}{p} + \frac{1}{q} = 1$ . Now let  $a$  and  $b$  be two positive numbers and consider  $f(t) = \frac{1}{p}(at)^p + \frac{1}{q}\left(\frac{b}{t}\right)^q$  for  $t > 0$ . Show the minimum value of  $f$  is  $ab$ . Prove the important inequality,  $ab \leq \frac{a^p}{p} + \frac{b^q}{q}$ .
21. Using Problem 20 establish the following magnificent inequality which is a case of Holder's inequality. For  $\frac{1}{p} + \frac{1}{q} = 1$ , and  $a_i, b_i$  positive numbers,

$$\sum_{i=1}^n a_i b_i \leq \left( \sum_{i=1}^n a_i^p \right)^{1/p} \left( \sum_{i=1}^n b_i^q \right)^{1/q}.$$

## 8.14 The Newton Raphson Method

The Newton Raphson method is a way to get approximations of solutions to various equations. For example, suppose you want to find  $\sqrt{2}$ . The existence of  $\sqrt{2}$  is not difficult to establish by considering the continuous function,  $f(x) = x^2 - 2$  which is negative at  $x = 0$  and positive at  $x = 2$ . Therefore, by the intermediate value theorem, there exists  $x \in (0, 2)$  such that  $f(x) = 0$  and this  $x$  must equal  $\sqrt{2}$ . The problem consists of how to find this number, not just to prove it exists. The following picture illustrates the procedure of the Newton Raphson method.



In this picture, a first approximation, denoted in the picture as  $x_1$  is chosen and then the tangent line to the curve  $y = f(x)$  at the point  $(x_1, f(x_1))$  is obtained. The equation of this tangent line is

$$y - f(x_1) = f'(x_1)(x - x_1).$$

Then extend this tangent line to find where it intersects the  $x$  axis. In other words, set  $y = 0$  and solve for  $x$ . This value of  $x$  is denoted by  $x_2$ . Thus

$$x_2 = x_1 - \frac{f(x_1)}{f'(x_1)}.$$

This second point,  $x_2$  is the second approximation and the same process is done for  $x_2$  that was done for  $x_1$  in order to get the third approximation,  $x_3$ . Thus

$$x_3 = x_2 - \frac{f(x_2)}{f'(x_2)}.$$

Continuing this way, yields a sequence of points,  $\{x_n\}$  given by

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}. \quad (8.24)$$

which hopefully has the property that  $\lim_{n \rightarrow \infty} x_n = x$  where  $f(x) = 0$ . You can see from the above picture that this must work out in the case of  $f(x) = x^2 - 2$ .

Now carry out the computations in the above case for  $x_1 = 2$  and  $f(x) = x^2 - 2$ . From (8.24),

$$x_2 = 2 - \frac{2}{4} = 1.5.$$

Then

$$x_3 = 1.5 - \frac{(1.5)^2 - 2}{2(1.5)} \leq 1.417,$$

$$x_4 = 1.417 - \frac{(1.417)^2 - 2}{2(1.417)} = 1.414216302046577,$$

What is the true value of  $\sqrt{2}$ ? To several decimal places this is  $\sqrt{2} = 1.414213562373095$ , showing that the Newton Raphson method has yielded a very good approximation after only a few iterations, even starting with an initial approximation, 2, which was not very good.

This method does not always work. For example, suppose you wanted to find the solution to  $f(x) = 0$  where  $f(x) = x^{1/3}$ . You should check that the sequence of iterates which results does not converge. This is because, starting with  $x_1$  the above procedure yields  $x_2 = -2x_1$  and so as the iteration continues, the sequence oscillates between positive and negative values as its absolute value gets larger and larger.



However, if  $f(x_0) = 0$  and  $f''(x) > 0$  for  $x$  near  $x_0$ , you can draw a picture to show that the method will yield a sequence which converges to  $x_0$  provided the first approximation,  $x_1$  is taken sufficiently close to  $x_0$ . Similarly, if  $f''(x) < 0$  for  $x$  near  $x_0$ , then the method produces a sequence which converges to  $x_0$  provided  $x_1$  is close enough to  $x_0$ .

## 8.15 Exercises

1. By drawing representative pictures, show convergence of the Newton Raphson method in the cases described above where  $f''(x) > 0$  near  $x_0$  or  $f''(x) < 0$  near  $x_0$ .
2. Draw some graphs to illustrate the Newton Raphson method does not yield a convergent sequence in the case where  $f(x) = x^{1/3}$ .
3. Use the Newton Raphson method to approximate the first positive solution of  $x - \tan x = 0$ . **Hint:** You may need to use a calculator to deal with  $\tan x$ .
4. Use the Newton Raphson method to compute an approximation to  $\sqrt{3}$  which is within  $10^{-6}$  of the true value. Explain how you know you are this close.
5. Using the Newton Raphson method and an appropriate picture, discuss the convergence of the recursively defined sequence  $x_{n+1} = ((p-1)x_n + cx_n^{1-p})/p$  where  $x_1, c > 0$  and  $p > 1$ .



# Antiderivatives And Differential Equations

A differential equation is an equation which involves an unknown function and its derivatives. Differential equations are the unifying idea in this chapter. Many interesting problems may be solved by formulating them as solutions of a suitable differential equation with initial condition called an initial value problem.

## 9.1 Initial Value Problems

The initial value problem is to find a function,  $y(x)$  for  $x \in [a, b]$  such that

$$y'(x) = f(x, y(x)), \quad y(a) = y_0.$$

Various assumptions are made on  $f(t, y)$ . At this time it is assumed that  $f$  does not depend on  $y$  and  $f$  is a given continuous defined  $[a, b]$ . Thus the initial value problem of interest here is one of the form

$$y'(x) = f(x), \quad y(a) = y_0. \quad (9.1)$$

**Theorem 9.1.1** *There is at most one solution to the initial value problem, (9.1) which is continuous on  $[a, b]$ .*

**Proof:** Suppose both  $A(x)$  and  $B(x)$  are solutions to this initial value problem for  $x \in (a, b)$ . Then letting  $H(x) \equiv A(x) - B(x)$ , it follows that  $H(a) = 0$  and  $H'(x) = A'(x) - B'(x) = f(x) - f(x) = 0$ . Therefore, from Corollary 6.8.4 on Page 133, it follows  $H(x)$  equals a constant on  $(a, b)$ . By continuity of  $H$ , this constant must equal  $H(a) = 0$ .

The main difficulty in solving these initial value problems like (9.1) is in finding a function whose derivative equals the given function. This is in general a very hard problem, although techniques for doing this are presented later which will cover many cases of interest. The functions whose derivatives equal a given function,  $f(x)$ , are called antiderivatives and there is a special notation used to denote them.

**Definition 9.1.2** *Let  $f$  be a function.  $\int f(x)dx$  denotes the set of antiderivatives of  $f$ . Thus  $F \in \int f(x)dx$  means  $F'(x) = f(x)$ . It is customary to refer to  $f(x)$  as the integrand. This symbol is also called the indefinite integral and sometimes is referred to as an integral although this last usage is not correct.*

The reason this last usage is not correct is that the integral of a function is a single number not a whole set of functions. Nevertheless, you can't escape the fact that it is common usage to call that symbol an integral.

**Lemma 9.1.3** Suppose  $F, G \in \int f(x) dx$  for  $x \in (a, b)$ . Then there exists a constant,  $C$  such that for all  $x \in (a, b)$ ,  $F(x) = G(x) + C$ .

**Proof:**

$$F'(x) - G'(x) = f(x) - f(x) = 0$$

for all  $x \in (a, b)$ . Consequently, by Corollary 6.8.4 on Page 133,  $F(x) - G(x) = C$ . This proves the lemma.

There is another simple lemma about antiderivatives.

**Lemma 9.1.4** If  $a$  and  $b$  are nonzero real numbers, and if  $\int f(x) dx$  and  $\int g(x) dx$  are nonempty, then

$$\int (af(x) + bg(x)) dx = a \int f(x) dx + b \int g(x) dx$$

**Proof:** The symbols on the two sides of the equation denote sets of functions. It is necessary to verify the two sets of functions are the same. Suppose then that  $F \in \int f(x) dx$  and  $G \in \int g(x) dx$ . Then  $aF + bG$  is a typical function of the right side of the equation. Taking the derivative of this function, yields  $af(x) + bg(x)$  and so this shows the set of functions on the right side is a subset of the set of functions on the left.

Now take  $H \in \int (af(x) + bg(x)) dx$  and pick  $F \in \int f(x) dx$ . Then  $aF \in a \int f(x) dx$  and

$$(H(x) - aF(x))' = af(x) + bg(x) - af(x) = bg(x)$$

showing that

$$H - aF \in \int bg(x) dx = b \int g(x) dx$$

because  $b \neq 0$ . Therefore

$$H \in aF + b \int g(x) dx \subseteq a \int f(x) dx + b \int g(x) dx.$$

This has shown the two sets of functions are the same and proves the lemma.

From Lemma 9.1.3 it follows that if  $F(x) \in \int f(x) dx$ , then every other function in  $\int f(x) dx$  is of the form  $F(x) + C$  for a suitable constant,  $C$ . Thus it is customary to write

$$\int f(x) dx = F(x) + C$$

where it is understood that  $C$  is an arbitrary constant, called a constant of integration.

From the formulas for derivatives presented earlier, the following table of antiderivatives follows.

$f(x)$	$\int f(x) dx$
$x^n, n \neq -1$	$\frac{x^{n+1}}{n+1} + C$
$x^{-1}$	$\ln x  + C$
$\cos(x)$	$\sin(x) + C$
$\sin(x)$	$-\cos(x) + C$
$e^x$	$e^x + C$
$\cosh(x)$	$\sinh(x) + C$
$\sinh(x)$	$\cosh(x) + C$
$\frac{1}{\sqrt{1-x^2}}$	$\arcsin(x) + C$
$\frac{1}{\sqrt{1+x^2}}$	$\operatorname{arcsinh} x + C$

The above table is a good starting point for other antiderivatives. For example,

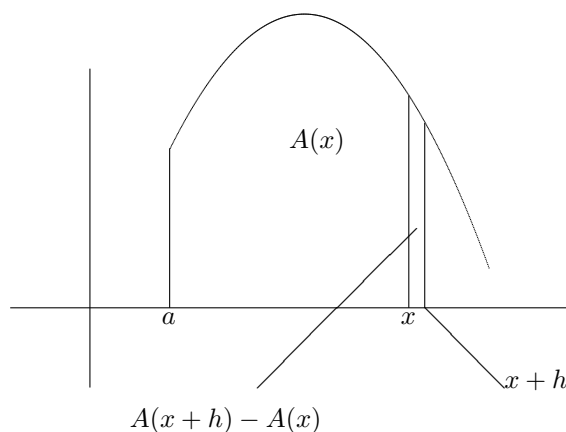
**Proposition 9.1.5** *Let  $\sum_{k=0}^n a_k x^k$  be a polynomial. Then*

$$\int \sum_{k=0}^n a_k x^k dx = \sum_{k=0}^n a_k \frac{x^{k+1}}{k+1} + C$$

**Proof:** This follows from the above table and Lemma 9.1.4.

## 9.2 Areas

Consider the problem of finding the area between the graph of a function of one variable and the  $x$  axis as illustrated in the following picture.



The curved line on the top represents the graph of the function  $y = f(x)$  and the symbol,  $A(x)$  represents the area between this curve and the  $x$  axis between the point,  $a$  and the point  $x$  as shown. The vertical line from the point,  $x+h$  up to the curve and the vertical line from  $x$  up to the curve define the area,  $A(x+h) - A(x)$  as indicated. You can see that this area is between  $hf(x)$  and  $hf(x+h)$ . This happens because the function is decreasing near  $x$ . In general, for continuous functions,  $f$ , Theorem 5.7.10 on Page 96 implies there exists  $x_M, x_m \in [x, x+h]$  with the properties

$$f(x_M) \equiv \max \{f(x) : x \in [x, x+h]\}$$

and

$$f(x_m) \equiv \min \{f(x) : x \in [x, x+h]\}$$

Then,

$$f(x_m) = \frac{hf(x_m)}{h} \leq \frac{A(x+h) - A(x)}{h} \leq \frac{hf(x_M)}{h} = f(x_M).$$

Therefore, using the squeezing theorem, Theorem 5.9.5, and the continuity of  $f$ ,

$$A'(x) \equiv \lim_{h \rightarrow 0} \frac{A(x+h) - A(x)}{h} = f(x).$$

The consideration of  $h < 0$  is also straightforward. This discussion implies the following theorem.

**Theorem 9.2.1** Let  $a < b$  and let  $f : [a, b] \rightarrow [0, \infty)$  be continuous. Then letting  $A(x)$  denote the area between  $a, x$ , the graph of the function, and the  $x$  axis,

$$A'(x) = f(x) \text{ for } x \in (a, b), \quad A(a) = 0. \quad (9.2)$$

Also,  $A$  is continuous on  $[a, b]$ .

The problem for  $A$  described in the above theorem is called an initial value problem and the equation,  $A'(x) = f(x)$  is a differential equation. It is called this because it is an equation for an unknown function,  $A(x)$  written in terms of the derivative of this unknown function. The assertion that  $A$  should be continuous on  $[a, b]$  follows from the fact that it has to be continuous on  $(a, b)$  because of the existence of its derivative and the above argument can also be used to obtain one sided derivatives for  $A$  at the end points,  $a$  and  $b$ , which yields continuity on  $[a, b]$ .

**Example 9.2.2** Let  $f(x) = x^2$  for  $x \in [1, 2]$ . Find the area between the graph of the function, the points 1 and 2, and the  $x$  axis.

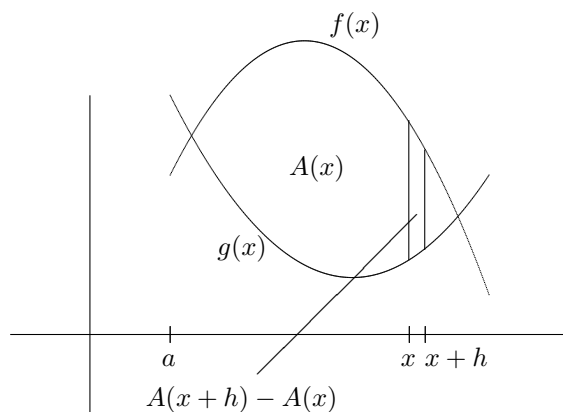
The function,  $\frac{x^3}{3} + C$  has the property that its derivative gives  $x^2$ . This is true for any  $C$ . It only remains to choose  $C$  in such a way that the function equals zero at  $x = 1$ . Thus  $C = -\frac{1}{3}$ . It follows that  $A(x) = \frac{x^3}{3} - \frac{1}{3}$ . Therefore, the area described equals  $A(2) = \frac{8}{3} - \frac{1}{3} = \frac{7}{3}$  square units.

**Example 9.2.3** Find the area between the graph of the function  $y = 1/x^2$  and the  $x$  axis for  $x$  between  $1/2$  and  $3$ .

The function  $-\frac{1}{x} + C$  has the property that its derivative equals  $1/x^2$ . Letting  $C = 2$ ,  $A(x) = -\frac{1}{x} + 2$  satisfies the appropriate initial value problem and so the area equals  $A(3) = \frac{5}{3}$ .

### 9.3 Area Between Graphs

It is a minor generalization to consider the area between the graphs of two functions. Consider the following picture.



You see that sometimes the function,  $f(x)$  is on top and sometimes the function,  $g(x)$  is on top. It is the length of the vertical line joining the two graphs which is of importance

and this length is always  $|f(x) - g(x)|$  regardless of which function is larger. By Theorem 5.7.10 on Page 96 there exist  $x_M, x_m \in [x, x+h]$  satisfying

$$|f(x_M) - g(x_M)| \equiv \max \{|f(x) - g(x)| : x \in [x, x+h]\}$$

and

$$|f(x_m) - g(x_m)| \equiv \min \{|f(x) - g(x)| : x \in [x, x+h]\}.$$

Then

$$\frac{|f(x_m) - g(x_m)|h}{h} \leq \frac{A(x+h) - A(x)}{h} \leq \frac{|f(x_M) - g(x_M)|h}{h},$$

and using the squeezing theorem, Theorem 5.9.5 on Page 100, as  $h \rightarrow 0$

$$A'(x) = |f(x) - g(x)|$$

Also  $A(a) = 0$  as before. This yields the following theorem which generalizes the one presented earlier because the  $x$  axis is the graph of the function,  $y = 0$ .

**Theorem 9.3.1** *Let  $a < b$  and let  $f, g : [a, b] \rightarrow \mathbb{R}$  be continuous. Then letting  $A(t)$  denote the area between the graphs of the two functions for  $x \in [a, t]$ ,*

$$A'(t) = |f(t) - g(t)| \text{ for } t \in (a, b), \quad A(a) = 0. \quad (9.3)$$

Also,  $A$  is continuous on  $[a, b]$ .

**Example 9.3.2** *Let  $f(x) = 8 - \frac{x^2}{2}$  and  $g(x) = \frac{x^2}{2} - 1$ . Find the area between the graphs of the two functions for  $x \in [-4, 3]$ .*

You should graph the two functions. The answer is  $A(3)$  where

$$A'(x) = |f(x) - g(x)| = |9 - x^2|, \quad A(-4) = 0.$$

Now

$$|9 - x^2| = \begin{cases} x^2 - 9 & \text{if } x \in [-4, -3] \\ 9 - x^2 & \text{if } x \in [-3, 3] \end{cases}$$

It follows that on  $(-4, -3)$ ,  $A(x) = \frac{x^3}{3} - 9x + C$  where  $C$  is chosen so that  $A(-4) = 0$ . Thus

$$\frac{(-4)^3}{3} - 9(-4) + C = 0$$

and so  $C = -\frac{44}{3}$ . Therefore, for  $x \in (-4, -3)$ ,

$$A(x) = \frac{x^3}{3} - 9x - \frac{44}{3}$$

Also,

$$A(-3) = \frac{(-3)^3}{3} - 9(-3) - \frac{44}{3} = \frac{10}{3}.$$

Now consider  $A(x)$  for  $x \in (-3, 3)$ . It is necessary to have  $A$  continuous on the interval  $[-4, 3]$  because it is supposed to have a derivative. Hence

$$A'(x) = 9 - x^2, \quad A(-3) = \frac{10}{3}$$

for  $x \in (-3, 3)$ . It follows that

$$A(x) = 9x - \frac{x^3}{3} + D$$

where

$$A(-3) = 9(-3) - \frac{(-3)^3}{3} + D = \frac{10}{3}$$

Thus  $D = \frac{64}{3}$  and  $A(x) = 9x - \frac{x^3}{3} + \frac{64}{3}$  for  $x \in (-3, 3)$ . By continuity of  $A$ , the area,  $A(3)$ , is given by

$$A(3) = 9(3) - \frac{(3)^3}{3} + \frac{64}{3} = \frac{118}{3}.$$

This illustrates the following procedure which you can use to find areas.

**Procedure 9.3.3** Suppose  $y = f(x)$  and  $y = g(x)$  are two functions defined for  $x \in [a, b]$ . Then to find the area between the two graphs of the functions, find an antiderivative of  $|f(x) - g(x)|$ ,  $A(x)$ . The desired area is then given by  $A(b) - A(a)$ .

**Proof:** Suppose  $B'(x) = |f(x) - g(x)|$ ,  $B(a) = 0$ . Then from the above explanation,  $B(b)$  is the desired area. From Lemma 9.1.3 there is some constant,  $C$  such that  $A(x) = B(x) + C$ . Therefore,

$$\begin{aligned} A(b) - A(a) &= (B(b) + C) - (B(a) + C) \\ &= (B(b) + C) - C = B(b) \end{aligned}$$

which is the desired area.

A similar procedure holds for finding the area between two functions which are of the form  $x = f(y)$  and  $x = g(y)$  for  $y \in [c, d]$ .

More generally, consider the initial value problem

$$F'(x) = f(x), F(a) = 0. \quad (9.4)$$

**Procedure 9.3.4** To solve the initial value problem, (9.4), find an antiderivative,  $G$ , such that  $G'(x) = f(x)$ . If there is a solution to the initial value problem, (9.4) it will equal  $G(x) - G(a)$ .

**Proof:** Suppose  $F$  solves (9.4). Then by Lemma 9.1.3, there exists a constant,  $C$  such that  $F(x) = G(x) + C$ . Since  $F(a) = 0$ , it follows  $C = -G(a)$  and so  $F(x) = G(x) - G(a)$ .

**Example 9.3.5** Find the area between  $x = 4 - y^2$  and  $x = -3y$ .

First find where the two graphs intersect.  $4 - y^2 = -3y$ . The solution is  $y = -1$  and  $4$ . For  $y$  in this interval, you can verify that  $4 - y^2 > -3y$  and so  $|4 - y^2 - (-3y)| = 4 - y^2 + 3y$ . An antiderivative is  $4y - \frac{y^3}{3} + \frac{3y^2}{2}$  and so the desired area is

$$4(4) - \frac{(4)^3}{3} + \frac{3(4)^2}{2} - \left( 4(-1) - \frac{(-1)^3}{3} + \frac{3(-1)^2}{2} \right) = \frac{125}{6}.$$

## 9.4 Exercises

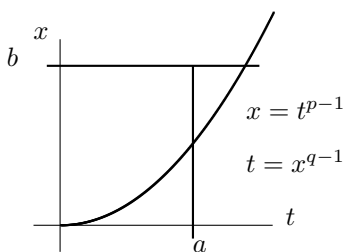
1. Find the area between the graphs of the functions,  $y = x^2 + 1$  and  $y = 3x + 5$ .
2. Find the area between the graphs of the functions,  $y = 2x^2$  and  $y = 6x + 8$ .
3. Find the area between the graphs of  $y = 5x + 14$  and  $y = x^2$ .
4. Find the area between the graphs of the functions  $y = x + 1$  and  $y = 2x$  for  $x \in [0, 3]$ .



5. Find the area between  $y = |x|$  and the  $x$  axis for  $x \in [-2, 2]$ .
6. Find the area between the graphs of  $y = x$  and  $y = \sin x$  for  $x \in [-\frac{\pi}{2}, \pi]$ .
7. Find the area between the graphs of  $x = y^2$  and  $y = 2 - x$ .
8. Find the area between the  $x$  axis and the graph of the function  $2x\sqrt{1+x^2}$  for  $x \in [0, 2]$ .  
**Hint:** Recall the chain rule for derivatives.
9. Show that the area of a right triangle equals one half the product of two sides which are not the hypotenuse.
10. Let  $A$  denote the region between the  $x$  axis and the graph of the function,  $f(x) = x - x^2$ . For  $k \in (0, 1)$ , the line  $y = kx$  divides this region into two pieces. Explain why there exists a number,  $k$  such that the area of these two pieces is exactly equal. **Hint:** This will likely involve the intermediate value theorem. Write an equation satisfied by  $k$  and then find an approximate value for  $k$ . **Hint:** You should draw plenty of pictures to do this last part.
11. Find the area between the graph of  $f(x) = 1/x$  for  $x \in [1, 2]$  and the  $x$  axis in terms of known functions.
12. Find the area between the graph of  $f(x) = 1/x^2$  for  $x \in [1, 2]$  and the  $x$  axis in terms of known functions.
13. Find the area between  $y = \sin x$  and  $y = \cos x$  for  $x \in [0, \frac{\pi}{4}]$ .
14. Find the area between  $e^x$  and  $\cos x$  for  $x \in [0, 2\pi]$ .
15. Find the area between  $y = e^x$  and  $y = 2x + 1$  for  $x > 0$ . In order to do this, you have to find a solution to  $e^x = 2x + 1$  and this will require a numerical procedure such as Newton's method or graphing and zooming on your calculator.
16. Find the area between  $y = \ln x$  and  $y = \sin x$  for  $x > 0$ . In order to do this, you have to find a solution to  $\ln x = \sin x$  and this will require a numerical procedure such as Newton's method or graphing and zooming on your calculator.
17. Let  $p > 1$ . An inequality which is of major importance is

$$ab \leq \frac{a^p}{p} + \frac{b^q}{q}$$

where here  $q$  is defined by  $1/p + 1/q = 1$ . Establish this inequality by adding up areas in the following picture.



In the picture the right side of the inequality represents the sum of all the areas and the left side is the area of the rectangle determined by  $(a, 0)$  and  $(0, b)$ .

## 9.5 The Method Of Substitution

Finding areas can often be reduced to solving an initial value problem. The crucial step in the solution of this initial value problem is to find an antiderivative for a given function. This is a hard problem if you insist on finding antiderivatives in terms of known functions but there are some general procedures for finding antiderivatives which are presented next. The method of substitution is based on the following formula which is merely a restatement of the chain rule.

$$\int f(g(x)) g'(x) dx = F(x) + C, \quad (9.5)$$

where  $F'(y) = f(y)$ . Here are some other examples of the method of substitution.

**Example 9.5.1** Find  $\int \sin(x) \cos(x) \sqrt{3 + 2^{-1} \sin^2(x)} dx$

Note it is of the form given in (9.5) with  $g(x) = 2^{-1} \sin^2(x)$  and  $F(u) = \frac{2}{3} \left( \sqrt{(3+u)} \right)^3$ . Therefore,

$$\int \sin(x) \cos(x) \sqrt{3 + 2^{-1} \sin^2(x)} dx = \frac{1}{12} \left( \sqrt{(12 + 2 \sin^2 x)} \right)^3 + C.$$

**Example 9.5.2** Find  $\int (1 + x^2)^6 2x dx$ .

This is a special case of (9.5) when  $g(x) = 1 + x^2$  and  $F(u) = u^7/7$ . Therefore, the answer is

$$\int (1 + x^2)^6 2x dx = (1 + x^2)^7 / 7 + C.$$

**Example 9.5.3** Find  $\int (1 + x^2)^6 x dx$

This equals

$$\frac{1}{2} \int (1 + x^2)^6 2x dx = \frac{1}{2} \frac{(1 + x^2)^7}{7} + C = \frac{(1 + x^2)^7}{14} + C.$$

Actually, it is not necessary to recall (9.5) and massage things to get them in that form. There is a trick based on the Leibniz notation for the derivative which is very useful and illustrated in the following example.

**Example 9.5.4** Find  $\int \cos(2x) \sin^2(2x) dx$ .

Let  $u = \sin(2x)$ . Then  $\frac{du}{dx} = 2 \cos(2x)$ . Now formally

$$\frac{du}{2} = \cos(2x) dx.$$

Thus

$$\begin{aligned} \int \cos(2x) \sin^2(2x) dx &= \frac{1}{2} \int u^2 du = \frac{u^3}{6} + C \\ &= \frac{(\sin(2x))^3}{6} + C \end{aligned}$$

The expression  $(1/2) du$  replaced the expression  $\cos(2x) dx$  which occurs in the original problem and the resulting problem in terms of  $u$  was much easier. This was solved and finally the original variable was replaced. When using this method, it is a good idea to check your answer to be sure you have not made a mistake. Thus in this example, the chain rule implies  $\left(\frac{(\sin(2x))^3}{6}\right)' = \cos(2x) \sin^2(2x)$  which verifies the answer is right. Here is another example.

**Example 9.5.5** Find  $\int \sqrt[3]{2x+7} x dx$ .

In this example  $u = 2x + 7$  so that  $du = 2dx$ . Then

$$\begin{aligned} \int \sqrt[3]{2x+7} x dx &= \int \sqrt[3]{u} \overbrace{\frac{u-7}{2}}^x \overbrace{\frac{1}{2}}^{dx} du \\ &= \int \left( \frac{1}{4} u^{4/3} - \frac{7}{4} u^{1/3} \right) du \\ &= \frac{3}{28} u^{7/3} - \frac{21}{16} u^{4/3} + C \\ &= \frac{3}{28} (2x+7)^{7/3} - \frac{21}{16} (2x+7)^{4/3} + C \end{aligned}$$

**Example 9.5.6** Find  $\int x 3^{x^2} dx$

Let  $u = 3^{x^2}$  so that  $\frac{du}{dx} = 2x \ln(3) 3^{x^2}$  and  $\frac{du}{2 \ln(3)} = x 3^{x^2} dx$ . Thus

$$\begin{aligned} \int x 3^{x^2} dx &= \frac{1}{2 \ln(3)} \int du = \frac{1}{2 \ln(3)} [u + C] \\ &= \frac{1}{2 \ln(3)} 3^{x^2} + \left( \frac{1}{2 \ln(3)} \right) C \end{aligned}$$

Since the constant is an arbitrary constant, this is written as

$$\frac{1}{2 \ln(3)} 3^{x^2} + C.$$

**Example 9.5.7** Find  $\int \cos^2(x) dx$

Recall that  $\cos(2x) = \cos^2(x) - \sin^2(x)$  and  $1 = \cos^2(x) + \sin^2(x)$ . Then subtracting and solving for  $\cos^2(x)$ ,

$$\cos^2(x) = \frac{1 + \cos(2x)}{2}.$$

Therefore,

$$\int \cos^2(x) dx = \int \frac{1 + \cos(2x)}{2} dx$$

Now letting  $u = 2x$ ,  $du = 2dx$  and so

$$\begin{aligned} \int \cos^2(x) dx &= \int \frac{1 + \cos(u)}{4} du \\ &= \frac{1}{4} u + \frac{1}{4} \sin u + C \\ &= \frac{1}{4} (2x + \sin(2x)) + C. \end{aligned}$$

Also

$$\int \sin^2(x) \, dx = -\frac{1}{2} \cos x \sin x + \frac{1}{2}x + C$$

which is left as an exercise.

**Example 9.5.8** Find  $\int \tan(x) \, dx$

Let  $u = \cos x$  so that  $du = -\sin(x) \, dx$ . Then writing the antiderivative in terms of  $u$ , this becomes  $\int \frac{-1}{u} \, du$ . At this point, recall that  $(\ln |u|)' = 1/u$ . Thus this antiderivative is  $-\ln |u| + C = \ln |u^{-1}| + C$  and so  $\int \tan(x) \, dx = \ln |\sec x| + C$ .

This illustrates a general procedure.

**Procedure 9.5.9**  $\int \frac{f'(x)}{f(x)} \, dx = \ln |f(x)| + C$ .

This follows from the chain rule.

**Example 9.5.10** Find  $\int \sec(x) \, dx$ .

This is usually done by a trick. You write as  $\int \frac{\sec(x)(\sec(x)+\tan(x))}{(\sec(x)+\tan(x))} \, dx$  and note that the numerator of the integrand is the derivative of the denominator. Thus  $\int \sec(x) \, dx = \ln |\sec(x) + \tan(x)| + C$ .

**Example 9.5.11** Find  $\int \csc(x) \, dx$ .

This is done like the antiderivatives for the secant.  $\frac{d}{dx} \csc(x) = -\csc(x) \cot(x)$  and  $\frac{d}{dx} \cot(x) = -\csc^2(x)$ . Write the integral as  $-\int \frac{\csc(x)(\cot(x)+\csc(x))}{(\cot(x)+\csc(x))} \, dx = -\ln |\cot(x) + \csc(x)| + C$ .

## 9.6 Exercises

1. Find the indicated antiderivatives.

- (a)  $\int \frac{x}{\sqrt{2x-3}} \, dx$
- (b)  $\int x(3x^2 + 6)^5 \, dx$
- (c)  $\int x \sin(x^2) \, dx$
- (d)  $\int \sin^3(2x) \cos(2x) \, dx$
- (e)  $\int \frac{1}{\sqrt{1+4x^2}} \, dx$  **Hint:** Remember the  $\sinh^{-1}$  function and its derivative.

2. Find the indicated antiderivatives.

- (a)  $\int \sec(3x) \, dx$
- (b)  $\int \sec^2(3x) \tan(3x) \, dx$
- (c)  $\int \frac{1}{3+5x^2} \, dx$
- (d)  $\int \frac{1}{\sqrt{5-4x^2}} \, dx$
- (e)  $\int \frac{3}{x\sqrt{4x^2-5}} \, dx$

3. Find the indicated antiderivatives.

- (a)  $\int x \cosh(x^2 + 1) \, dx$

- (b)  $\int x^3 5^{x^4} dx$   
 (c)  $\int \sin(x) 7^{\cos(x)} dx$   
 (d)  $\int x \sin(x^2) dx$   
 (e)  $\int x^5 \sqrt{2x^2 + 1} dx$  **Hint:** Let  $u = 2x^2 + 1$ .
4. Find  $\int \sin^2(x) dx$ . **Hint:** Derive and use  $\sin^2(x) = \frac{1 - \cos(2x)}{2}$ .
5. Find the area between the graphs of  $y = \sin(2x)$  and  $y = \cos(2x)$  for  $x \in [0, 2\pi]$ .
6. Find the area between the graphs of  $y = \sin^2 x$ , and the  $x$  axis, for  $x \in [0, 2\pi]$ .
7. Find the area between the graphs of  $y = \cos^2 x$  and  $y = \sin^2 x$  for  $x \in [0, \pi/2]$ .
8. Find the indicated antiderivatives.
- (a)  $\int \frac{\ln x}{x} dx$   
 (b)  $\int \frac{x^3}{3+x^4} dx$   
 (c)  $\int \frac{1}{x^2+2x+2} dx$  **Hint:** Complete the square in the denominator and then let  $u = x + 1$ .  
 (d)  $\int \frac{1}{\sqrt{4-x^2}} dx$   
 (e)  $\int \frac{1}{x\sqrt{x^2-9}} dx$  **Hint:** Let  $x = 3u$ .  
 (f)  $\int \frac{\ln(x^2)}{x} dx$   
 (g) Find  $\int \frac{x^3}{\sqrt{(6x^2+5)}} dx$   
 (h) Find  $\int x^3 \sqrt[3]{(6x+4)} dx$
9. Find the indicated antiderivatives.
- (a)  $\int x\sqrt{(2x+4)} dx$   
 (b)  $\int x\sqrt{(3x+2)} dx$   
 (c)  $\int \frac{1}{\sqrt{(36-25x^2)}} dx$   
 (d)  $\int \frac{1}{x} \sqrt{(3x+5)} dx$   
 (e)  $\int \frac{1}{\sqrt{(9-4x^2)}} dx$   
 (f)  $\int \frac{1}{\sqrt{(1+4x^2)}} dx$   
 (g)  $\int 4 \frac{x}{\sqrt{(3x-1)}} dx$   
 (h)  $\int \frac{1}{x^2} \sqrt[3]{(6x+4)} dx$   
 (i)  $\int 4 \frac{x}{\sqrt{(5x+1)}} dx$   
 (j)  $\int \frac{1}{x\sqrt{(9x^2-4)}} dx$   
 (k)  $\int \frac{1}{2\sqrt{(9+4x^2)}} dx$
10. Find the area between the graph of  $f(x) = x^3 / (2 + 3x^4)$  and the  $x$  axis for  $x \in [0, 4]$ .
11. Find  $\int \frac{1}{x^{1/3} + x^{1/2}} dx$ . **Hint:** Try letting  $x = u^6$  and use long division.

## 9.7 Integration By Parts

Another technique for finding antiderivatives is called integration by parts and is based on the product rule. Recall the product rule. If  $u'$  and  $v'$  exist, then

$$(uv)'(x) = u'(x)v(x) + u(x)v'(x). \quad (9.6)$$

Therefore,

$$(uv)'(x) - u'(x)v(x) = u(x)v'(x)$$

**Proposition 9.7.1** *Let  $u$  and  $v$  be differentiable functions for which  $\int u(x)v'(x) dx$  and  $\int u'(x)v(x) dx$  are nonempty. Then*

$$uv - \int u'(x)v(x) dx = \int u(x)v'(x) dx. \quad (9.7)$$

**Proof:** Let  $F \in \int u'(x)v(x) dx$ . Then

$$(uv - F)' = (uv)' - F' = (uv)' - u'v = uv'$$

by the chain rule. Therefore every function from the left in (9.7) is a function found in the right side of (9.7). Now let  $G \in \int u(x)v'(x) dx$ . Then  $(uv - G)' = -uv' + (uv)' = u'v$  by the chain rule. It follows that  $uv - G \in \int u'(x)v(x) dx$  and so  $G \in uv - \int u'(x)v(x) dx$ . Thus every function from the right in (9.7) is a function from the left. This proves the proposition.

**Example 9.7.2** Find  $\int x \sin(x) dx$

Let  $u(x) = x$  and  $v'(x) = \sin(x)$ . Then applying (9.7),

$$\begin{aligned} \int x \sin(x) dx &= (-\cos(x))x - \int (-\cos(x)) dx \\ &= -x \cos(x) + \sin(x) + C. \end{aligned}$$

**Example 9.7.3** Find  $\int x \ln(x) dx$

Let  $u(x) = \ln(x)$  and  $v'(x) = x$ . Then from (9.7),

$$\begin{aligned} \int x \ln(x) dx &= \frac{x^2}{2} \ln(x) - \int \frac{x^2}{2} \left( \frac{1}{x} \right) \\ &= \frac{x^2}{2} \ln(x) - \int \frac{x}{2} \\ &= \frac{x^2}{2} \ln(x) - \frac{1}{4}x^2 + C \end{aligned}$$

**Example 9.7.4** Find  $\int \arctan(x) dx$

Let  $u(x) = \arctan(x)$  and  $v'(x) = 1$ . Then from (9.7),

$$\begin{aligned} \int \arctan(x) dx &= x \arctan(x) - \int x \left( \frac{1}{1+x^2} \right) dx \\ &= x \arctan(x) - \frac{1}{2} \int \frac{2x}{1+x^2} dx \\ &= x \arctan(x) - \frac{1}{2} \ln(1+x^2) + C. \end{aligned}$$

## 9.8 Exercises

1. Find the following antiderivatives.

(a)  $\int x e^{-3x} dx$

(b)  $\int \frac{1}{x(\ln(|x|))^2} dx$

(c)  $\int x\sqrt{2-x} dx$

(d)  $\int (\ln|x|)^2 dx$  **Hint:** Let  $u(x) = (\ln|x|)^2$  and  $v'(x) = 1$ .

(e)  $\int x^3 \cos(x^2) dx$

2. Show that  $\int \sec^3(x) dx = \frac{1}{2} \tan(x) \sec(x) + \frac{1}{2} \ln|\sec x + \tan x| + C$ .

3. Consider the following argument. Integrate by parts, letting  $u(x) = x$  and  $v'(x) = \frac{1}{x^2}$  to get

$$\begin{aligned} \int \frac{1}{x} dx &= \int x \left( \frac{1}{x^2} \right) dx = \left( -\frac{1}{x} \right) x + \int \frac{1}{x} dx \\ &= -1 + \int \frac{1}{x} dx. \end{aligned}$$

Now subtracting  $\int \frac{1}{x} dx$  from both sides,  $0 = -1$ . Is there anything wrong here? If so, what?

4. Find the following antiderivatives.

(a)  $\int x^3 \arctan(x) dx$

(b)  $\int x^3 \ln(x) dx$

(c)  $\int x^2 \sin(x) dx$

(d)  $\int x^2 \cos(x) dx$

(e)  $\int x \arcsin(x) dx$

(f)  $\int \cos(2x) \sin(3x) dx$

(g)  $\int x^3 e^{x^2} dx$

(h)  $\int x^3 \cos(x^2) dx$

5. Find the antiderivatives

(a)  $\int x^2 \sin x dx$

(b)  $\int x^2 \sin x dx$

(c)  $\int x^3 7^x dx$

(d)  $\int x^2 \ln x dx$

(e)  $\int (x+2)^2 e^x dx$

(f)  $\int x^3 2^x dx$

(g)  $\int \sec^3(2x) \tan(2x) dx$

(h)  $\int x^2 7^x dx$

6. Try doing  $\int \sin^2 x dx$  the obvious way. If you don't make any mistakes, the process will go in circles. Now do it by taking  $\int \sin^2 x dx = x \sin^2 x - 2 \int x \sin x \cos x dx = x \sin^2 x - \int x \sin(2x) dx$ .

## 9.9 Trig. Substitutions

Certain antiderivatives are easily obtained by making an auspicious substitution involving a trig. function. The technique will be illustrated by presenting examples.

**Example 9.9.1** Find  $\int \frac{1}{(x^2+2x+2)^2} dx$ .

Complete the square as before and write

$$\int \frac{1}{(x^2 + 2x + 2)^2} dx = \int \frac{1}{((x+1)^2 + 1)^2} dx$$

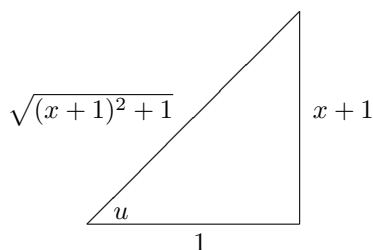
Use the following substitution next.

$$x + 1 = \tan u \tag{9.8}$$

so  $dx = (\sec^2 u) du$ . Therefore, this last indefinite integral becomes

$$\begin{aligned} \int \frac{\sec^2 u}{(\tan^2 u + 1)^2} du &= \int (\cos^2 u) du \\ &= \int \frac{1 + \cos 2u}{2} du \\ &= \frac{u}{2} + \frac{\sin 2u}{4} + C \\ &= \frac{u}{2} + \frac{2 \sin u \cos u}{4} + C \end{aligned}$$

Next write this in terms of  $x$  using the following device based on the following picture.



In this picture which is descriptive of (9.8),  $\sin u = \frac{x+1}{\sqrt{(x+1)^2 + 1}}$  and  $\cos u = \frac{1}{\sqrt{(x+1)^2 + 1}}$ . Therefore, putting in this information to change back to the  $x$  variable,

$$\begin{aligned} &\int \frac{1}{(x^2 + 2x + 2)^2} dx \\ &= \frac{1}{2} \arctan(x+1) + \frac{1}{2} \frac{x+1}{\sqrt{(x+1)^2 + 1}} \frac{1}{\sqrt{(x+1)^2 + 1}} + C \\ &= \frac{1}{2} \arctan(x+1) + \frac{1}{2} \frac{x+1}{(x+1)^2 + 1} + C. \end{aligned}$$

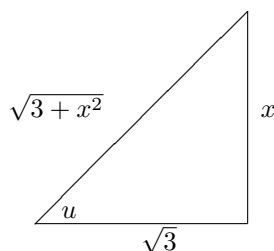


**Example 9.9.2** Find  $\int \frac{1}{\sqrt{x^2+3}} dx$ .

Let  $x = \sqrt{3} \tan u$  so  $dx = \sqrt{3} (\sec^2 u) du$ . Making the substitution, consider

$$\begin{aligned} & \int \frac{1}{\sqrt{3}\sqrt{\tan^2 u + 1}} \sqrt{3} (\sec^2 u) du \\ &= \int (\sec u) du = \ln |\sec u + \tan u| + C \end{aligned}$$

Now the following diagram is descriptive of the above transformation.



Using the above diagram,  $\sec u = \frac{\sqrt{3+x^2}}{\sqrt{3}}$  and  $\tan u = \frac{x}{\sqrt{3}}$ . Therefore, restoring the  $x$  variable,

$$\begin{aligned} \int \frac{1}{\sqrt{x^2+3}} dx &= \ln \left| \frac{\sqrt{3+x^2}}{\sqrt{3}} + \frac{x}{\sqrt{3}} \right| + C \\ &= \ln \left| \sqrt{3+x^2} + x \right| + C. \end{aligned}$$

**Example 9.9.3** Find  $\int (4x^2 + 3)^{1/2} dx$ .

Let  $2x = \sqrt{3} \tan u$  so  $2dx = \sqrt{3} \sec^2(u) du$ . Then making the substitution,

$$\begin{aligned} & \sqrt{3} \int (\tan^2 u + 1)^{1/2} \frac{\sqrt{3}}{2} \sec^2(u) du \\ &= \frac{3}{2} \int \sec^3(u) du. \end{aligned} \tag{9.9}$$

Now use integration by parts to obtain

$$\begin{aligned} \int \sec^3(u) du &= \int \sec^2(u) \sec(u) du = \\ &= \tan(u) \sec(u) - \int \tan^2(u) \sec(u) du \\ &= \tan(u) \sec(u) - \int (\sec^2(u) - 1) \sec(u) du \\ &= \tan(u) \sec(u) + \int \sec(u) du - \int \sec^3(u) du \\ &= \tan(u) \sec(u) + \ln |\sec(u) + \tan(u)| - \int \sec^3(u) du \end{aligned}$$

Therefore,

$$2 \int \sec^3(u) \, du = \tan(u) \sec(u) + \ln |\sec(u) + \tan(u)| + C$$

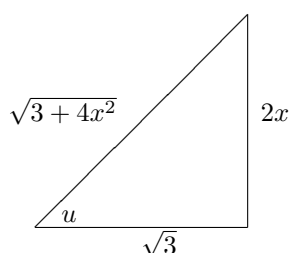
and so

$$\int \sec^3(u) \, du = \frac{1}{2} [\tan(u) \sec(u) + \ln |\sec(u) + \tan(u)|] + C. \quad (9.10)$$

Now it follows from (9.9) that in terms of  $u$  the set of antiderivatives is given by

$$\frac{3}{4} [\tan(u) \sec(u) + \ln |\sec(u) + \tan(u)|] + C$$

Use the following diagram to change back to the variable,  $x$ .



From the diagram,  $\tan(u) = \frac{2x}{\sqrt{3}}$  and  $\sec(u) = \frac{\sqrt{3+4x^2}}{\sqrt{3}}$ . Therefore,

$$\begin{aligned} \int (4x^2 + 3)^{1/2} \, dx &= \\ &= \frac{3}{4} \left[ \frac{2x}{\sqrt{3}} \frac{\sqrt{3+4x^2}}{\sqrt{3}} + \ln \left| \frac{\sqrt{3+4x^2}}{\sqrt{3}} + \frac{2x}{\sqrt{3}} \right| \right] + C \\ &= \frac{3}{4} \left[ \frac{2x}{\sqrt{3}} \frac{\sqrt{3+4x^2}}{\sqrt{3}} + \ln \left| \frac{\sqrt{3+4x^2}}{\sqrt{3}} + \frac{2x}{\sqrt{3}} \right| \right] + C \\ &= \frac{1}{2} x \sqrt{3+4x^2} + \frac{3}{4} \ln |\sqrt{3+4x^2} + 2x| + C \end{aligned}$$

Note that these examples involved something of the form  $(a^2 + (bx)^2)$  and the trig substitution,

$$bx = a \tan u$$

was the right one to use. This is the auspicious substitution which often simplifies these sorts of problems.

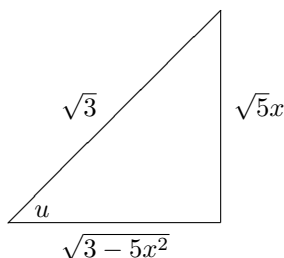
**Example 9.9.4** Find  $\int \sqrt{3 - 5x^2} \, dx$

In this example, let  $\sqrt{5}x = \sqrt{3} \sin(u)$  so  $\sqrt{5}dx = \sqrt{3} \cos(u) \, du$ . The reason this might be a good idea is that it will get rid of the square root sign as shown below. Making the

substitution,

$$\begin{aligned}
 \int \sqrt{3-5x^2} \, dx &= \sqrt{3} \int \sqrt{1-\sin^2(u)} \frac{\sqrt{3}}{\sqrt{5}} \cos(u) \, du \\
 &= \frac{3}{\sqrt{5}} \int \cos^2(u) \, du \\
 &= \frac{3}{\sqrt{5}} \int \frac{1+\cos 2u}{2} \, du \\
 &= \frac{3}{\sqrt{5}} \left( \frac{u}{2} + \frac{\sin 2u}{4} \right) + C \\
 &= \frac{3}{2\sqrt{5}} u + \frac{3}{2\sqrt{5}} \sin u \cos u + C
 \end{aligned}$$

The appropriate diagram is the following.



From the diagram,  $\sin(u) = \frac{\sqrt{5}x}{\sqrt{3}}$  and  $\cos(u) = \frac{\sqrt{3-5x^2}}{\sqrt{3}}$ . Therefore, changing back to  $x$ ,

$$\begin{aligned}
 \int \sqrt{3-5x^2} \, dx &= \\
 &= \frac{3}{2\sqrt{5}} \arcsin\left(\frac{\sqrt{5}x}{\sqrt{3}}\right) + \frac{3}{2\sqrt{5}} \frac{\sqrt{5}x}{\sqrt{3}} \frac{\sqrt{3-5x^2}}{\sqrt{3}} + C \\
 &= \frac{3}{10} \sqrt{5} \arcsin\left(\frac{1}{3} \sqrt{15}x\right) + \frac{1}{2} x \sqrt{(3-5x^2)} + C
 \end{aligned}$$

**Example 9.9.5** Find  $\int \sqrt{5x^2-3} \, dx$

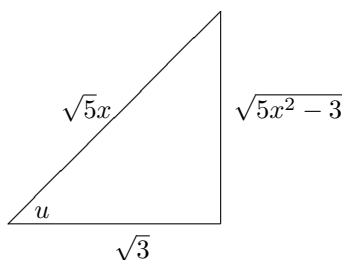
In this example, let  $\sqrt{5}x = \sqrt{3} \sec(u)$  so  $\sqrt{5}dx = \sqrt{3} \sec(u) \tan(u) \, du$ . Then changing the variable, consider

$$\begin{aligned}
 &\sqrt{3} \int \sqrt{\sec^2(u)-1} \frac{\sqrt{3}}{\sqrt{5}} \sec(u) \tan(u) \, du \\
 &= \frac{3}{\sqrt{5}} \int \tan^2(u) \sec(u) \, du \\
 &= \frac{3}{\sqrt{5}} \left[ \int \sec^3(u) \, du - \int \sec(u) \, du \right].
 \end{aligned}$$

Now from (9.10), this equals

$$\begin{aligned}
 &\frac{3}{\sqrt{5}} \left[ \frac{1}{2} [\tan(u) \sec(u) + \ln |\sec(u) + \tan(u)|] - \ln |\tan(u) + \sec(u)| \right] + C \\
 &= \frac{3}{2\sqrt{5}} \tan(u) \sec(u) - \frac{3}{2\sqrt{5}} \ln |\sec(u) + \tan(u)| + C.
 \end{aligned}$$

Now it is necessary to change back to  $x$ . The diagram is as follows.



Therefore,  $\tan(u) = \frac{\sqrt{5x^2-3}}{\sqrt{3}}$  and  $\sec(u) = \frac{\sqrt{5x}}{\sqrt{3}}$  and so

$$\begin{aligned} \int \sqrt{5x^2-3} \, dx &= \\ &= \frac{3}{2\sqrt{5}} \frac{\sqrt{5x^2-3}}{\sqrt{3}} \frac{\sqrt{5x}}{\sqrt{3}} - \frac{3}{2\sqrt{5}} \ln \left| \frac{\sqrt{5x}}{\sqrt{3}} + \frac{\sqrt{5x^2-3}}{\sqrt{3}} \right| + C \\ &= \frac{1}{2} \left( \sqrt{5x^2-3} \right) x - \frac{3}{10} \sqrt{5} \ln \left| \sqrt{5x} + \sqrt{(-3+5x^2)} \right| + C \end{aligned}$$

To summarize, here is a short table of auspicious substitutions corresponding to certain expressions.

Expression	$a^2 + b^2 x^2$	$a^2 - b^2 x^2$	$a^2 x^2 - b^2$
Trig. substitution	$bx = a \tan(u)$	$bx = a \sin(u)$	$ax = b \sec(u)$

Of course there are no “magic bullets” but these substitutions will often simplify an expression enough to allow you to find an antiderivative. These substitutions are often especially useful when the expression is enclosed in a square root.

## 9.10 Exercises

1. Find the antiderivatives.

- (a)  $\int \frac{x}{\sqrt{(4-x^2)}} \, dx$
- (b)  $\int \frac{3}{\sqrt{(36-25x^2)}} \, dx$
- (c)  $\int \frac{3}{\sqrt{(16-25x^2)}} \, dx$
- (d)  $\int \frac{1}{\sqrt{(4-9x^2)}} \, dx$
- (e)  $\int \frac{1}{\sqrt{(36-x^2)}} \, dx$
- (f)  $\int \left( \sqrt{(9-16x^2)} \right)^3 \, dx$
- (g)  $\int \left( \sqrt{(16-x^2)} \right)^5 \, dx$
- (h)  $\int \sqrt{(25-36x^2)} \, dx$

$$(i) \int \left( \sqrt{4-9x^2} \right)^3 dx$$

$$(j) \int \sqrt{1-9x^2} dx$$

2. Find the antiderivatives.

$$(a) \int \sqrt{36x^2 - 25} dx$$

$$(b) \int \sqrt{x^2 - 4} dx$$

$$(c) \int \left( \sqrt{16x^2 - 9} \right)^3 dx$$

$$(d) \int \sqrt{25x^2 - 16} dx$$

3. Find the antiderivatives.

$$(a) \int \frac{1}{26+x^2-2x} dx \quad \textbf{Hint:} \text{ Complete the square.}$$

$$(b) \int \sqrt{x^2 + 9} dx$$

$$(c) \int \sqrt{4x^2 + 25} dx$$

$$(d) \int x\sqrt{4x^4 + 9} dx$$

$$(e) \int x^3\sqrt{4x^4 + 9} dx$$

$$(f) \int \frac{1}{(25+36(2x-3)^2)^2} dx$$

$$(g) \int \frac{1}{(16+25(x-3)^2)^2} dx$$

$$(h) \int \frac{1}{261+25x^2-150x} dx \quad \textbf{Hint:} \text{ Complete the square.}$$

$$(i) \int \left( \sqrt{25x^2 + 9} \right)^3 dx$$

$$(j) \int \frac{1}{25+16x^2} dx$$

4. Find the antiderivatives. **Hint:** Complete the square.

$$(a) \int \sqrt{4x^2 + 16x + 15} dx$$

$$(b) \int \sqrt{x^2 + 6x} dx$$

$$(c) \int \frac{3}{\sqrt{(-32-9x^2-36x)}} dx$$

$$(d) \int \frac{3}{\sqrt{(-5-x^2-6x)}} dx$$

$$(e) \int \frac{1}{\sqrt{(9-16x^2-32x)}} dx$$

$$(f) \int \sqrt{4x^2 + 16x + 7} dx$$

## 9.11 Partial Fractions

The main technique for finding antiderivatives in the case  $f(x) = \frac{p(x)}{q(x)}$  for  $p$  and  $q$  polynomials is the technique of partial fractions. Before presenting this technique, a few more examples are presented.

**Example 9.11.1** Find  $\int \frac{1}{x^2+2x+2} dx$ .

To do this complete the square in the denominator to write

$$\int \frac{1}{x^2 + 2x + 2} dx = \int \frac{1}{(x+1)^2 + 1} dx$$

Now change the variable, letting  $u = x + 1$  so that  $du = dx$ . Then the last indefinite integral reduces to

$$\int \frac{1}{u^2 + 1} du = \arctan u + C$$

and so

$$\int \frac{1}{x^2 + 2x + 2} dx = \arctan(x + 1) + C.$$

**Example 9.11.2** Find  $\int \frac{1}{3x+5} dx$ .

Let  $u = 3x + 5$  so  $du = 3dx$  and changing the variable,

$$\frac{1}{3} \int \frac{1}{u} du = \frac{1}{3} \ln |u| + C.$$

Therefore,

$$\int \frac{1}{3x+5} dx = \frac{1}{3} \ln |3x+5| + C.$$

**Example 9.11.3** Find  $\int \frac{3x+2}{x^2+x+1} dx$ .

First complete the square in the denominator.

$$\begin{aligned} \int \frac{3x+2}{x^2+x+1} dx &= \int \frac{3x+2}{x^2+x+\frac{1}{4}+\frac{3}{4}} dx \\ &= \int \frac{3x+2}{\left(x+\frac{1}{2}\right)^2+\frac{3}{4}} dx. \end{aligned}$$

Now let

$$\left(x + \frac{1}{2}\right)^2 = \frac{3}{4}u^2$$

so that  $x + \frac{1}{2} = \frac{\sqrt{3}}{2}u$ . Therefore,  $dx = \frac{\sqrt{3}}{2}du$  and changing the variable,

$$\begin{aligned} &\frac{4}{3} \int \frac{3\left(\frac{\sqrt{3}}{2}u - \frac{1}{2}\right) + 2}{u^2 + 1} \frac{\sqrt{3}}{2} du \\ &= \frac{\sqrt{3}}{2} \left( 2\sqrt{3} \int \frac{u}{u^2 + 1} du - \frac{2}{3} \int \frac{1}{u^2 + 1} du \right) \\ &= \frac{\sqrt{3}}{2} \left( \sqrt{3} \int \frac{2u}{u^2 + 1} du - \frac{2}{3} \int \frac{1}{u^2 + 1} du \right) \\ &= \frac{3}{2} \ln(u^2 + 1) - \frac{\sqrt{3}}{3} \arctan u + C \end{aligned}$$

Therefore,

$$\begin{aligned} &\int \frac{3x+2}{x^2+x+1} dx = \\ &\frac{3}{2} \ln \left( \left( \frac{2}{\sqrt{3}} \left( x + \frac{1}{2} \right) \right)^2 + 1 \right) - \frac{\sqrt{3}}{3} \arctan \left( \frac{2}{\sqrt{3}} \left( x + \frac{1}{2} \right) \right) + C. \end{aligned}$$

The following simple but important Lemma is needed to continue.

**Lemma 9.11.4** *Let  $f(x)$  and  $g(x)$  be polynomials. Then there exists a polynomial,  $q(x)$  such that*

$$f(x) = q(x)g(x) + r(x)$$

where the degree of  $r(x) < \text{degree of } g(x)$  or  $r(x) = 0$ .

**Proof:** Consider the polynomials of the form  $f(x) - g(x)l(x)$  and out of all these polynomials, pick the one which has the smallest degree. This can be done because of the well ordering of the natural numbers discussed earlier. Let this take place when  $l(x) = q_1(x)$  and let  $r(x) = f(x) - g(x)q_1(x)$ . It is required to show degree of  $r(x) < \text{degree of } g(x)$  or else  $r(x) = 0$ . Suppose  $f(x) - g(x)l(x)$  is never equal to zero for any  $l(x)$ . Then  $r(x) \neq 0$ . It is required to show the degree of  $r(x)$  is smaller than the degree of  $g(x)$ . If this doesn't happen, then the degree of  $r \geq$  the degree of  $g$ . Let

$$\begin{aligned} r(x) &= b_mx^m + \cdots + b_1x + b_0 \\ g(x) &= a_nx^n + \cdots + a_1x + a_0 \end{aligned}$$

where  $m \geq n$  and  $b_m$  and  $a_n$  are nonzero. Then letting

$$\begin{aligned} r_1(x) &= r(x) - \frac{x^{m-n}b_m}{a_n}g(x) \\ &= f(x) - g(x)q_1(x) - \frac{x^{m-n}b_m}{a_n}g(x) \\ &= f(x) - g(x)\left(q_1(x) + \frac{x^{m-n}b_m}{a_n}\right), \end{aligned}$$

it follows this is not zero by the assumption that  $f(x) - g(x)l(x)$  is never equal to zero for any  $l(x)$ . Now the degree of  $r_1(x) < \text{degree of } r(x)$ , a contradiction to the construction of  $r(x)$ . This proves the lemma.

**Corollary 9.11.5** *Let  $f(x)$  and  $g(x)$  be polynomials. Then there exists a polynomial,  $r(x)$  such that the degree of  $r(x) < \text{degree of } g(x)$  and a polynomial,  $q(x)$  such that*

$$\frac{f(x)}{g(x)} = q(x) + \frac{r(x)}{g(x)}.$$

**Example 9.11.6** Find  $\int \frac{-x^3+11x^2+24x+14}{(2x+3)(x+5)(x^2+x+1)} dx$ .

In this problem, first check to see if the degree of the numerator in the integrand is less than the degree of the denominator. In this case, this is so. If it is not so, use long division to write the integrand as the sum of a polynomial with a rational function in which the degree of the numerator is less than the degree of the denominator. See the preceding corollary which guarantees this can be done. Now look for a partial fractions expansion for the integrand which is in the following form.

$$\frac{a}{2x+3} + \frac{b}{x+5} + \frac{cx+d}{x^2+x+1}$$

and try to find constants,  $a, b, c$ , and  $d$  so that the above rational functions sum to the integrand. The reason  $cx+d$  is used in the numerator of the last expression is that  $x^2+x+1$  cannot be factored using real polynomials. Thus the problem involves finding  $a, b, c, d$ , such that

$$\frac{-x^3+11x^2+24x+14}{(2x+3)(x+5)(x^2+x+1)} = \frac{a}{2x+3} + \frac{b}{x+5} + \frac{cx+d}{x^2+x+1}$$

and so

$$\begin{aligned} -x^3 + 11x^2 + 24x + 14 &= a(x+5)(x^2+x+1) + \\ &+ b(2x+3)(x^2+x+1) + (cx+d)(2x+3)(x+5). \end{aligned} \quad (9.11)$$

Now these are two polynomials which are supposed to be equal. Therefore, they have the same coefficients. Multiplying the right side out and collecting the terms,

$$\begin{aligned} -x^3 + 11x^2 + 24x + 14 &= \\ &= (2b+2c+a)x^3 + (6a+5b+13c+2d)x^2 + (6a+13d+5b+15c)x + 15d+5a+3b \end{aligned}$$

and therefore, it is necessary to solve the equations,

$$\begin{aligned} 2b+2c+a &= -1 \\ 6a+5b+13c+2d &= 11 \\ 6a+13d+5b+15c &= 24 \\ 15d+5a+3b &= 14 \end{aligned}$$

The solution is  $c=1, a=1, b=-2, d=1$ . Therefore,

$$\frac{-x^3 + 11x^2 + 24x + 14}{(2x+3)(x+5)(x^2+x+1)} = \frac{1}{2x+3} - \frac{2}{x+5} + \frac{1+x}{x^2+x+1}.$$

This may look like a fairly formidable problem. In reality it is not that bad. First let  $x = -5$  in (9.11) and obtain a simple equation for finding  $b$ . Next let  $x = -3/2$  to get a simple equation for  $a$ . This reduces the above system to a more manageable size. Anyway, it is now possible to find the antiderivative of the given function.

$$\begin{aligned} \int \frac{-x^3 + 11x^2 + 24x + 14}{(2x+3)(x+5)(x^2+x+1)} dx &= \\ \int \frac{1}{2x+3} dx - \int \frac{2}{x+5} dx + \int \frac{1+x}{x^2+x+1} dx. \end{aligned}$$

Now each of these indefinite integrals can be found using the techniques given above. Thus the desired antiderivative equals

$$\begin{aligned} &\frac{1}{2} \ln |2x+3| - 2 \ln |x+5| + \\ &+ \frac{1}{2} \ln (x^2+x+1) + \frac{1}{3} \sqrt{3} \arctan \left( \frac{\sqrt{3}}{3} (2x+1) \right) + C. \end{aligned}$$

This was a long example. Here is an easier one.

**Example 9.11.7** Find  $\int \frac{3x^5+7}{x^2-1} dx$ .

In this case the degree of the numerator is larger than the degree of the denominator and so long division must first be used. Thus

$$\frac{3x^5+7}{x^2-1} = 3x^3 + 3x + \frac{7+3x}{x^2-1}$$

Now look for a partial fractions expansion of the form

$$\frac{7+3x}{x^2-1} = \frac{a}{(x-1)} + \frac{b}{(x+1)}.$$



Therefore,

$$7 + 3x = a(x + 1) + b(x - 1).$$

Letting  $x = 1$ ,  $a = 5$ . Then letting  $x = -1$ , it follows  $b = -2$ . Therefore,

$$\frac{7 + 3x}{x^2 - 1} = \frac{5}{x - 1} - \frac{2}{x + 1}$$

and so

$$\frac{3x^5 + 7}{x^2 - 1} = 3x^3 + 3x + \frac{5}{x - 1} - \frac{2}{x + 1}.$$

therefore,

$$\int \frac{3x^5 + 7}{x^2 - 1} dx = \frac{3}{4}x^4 + \frac{3}{2}x^2 + 5 \ln(x - 1) - 2 \ln(x + 1) + C.$$

What is done when the factors are repeated?

**Example 9.11.8** Find  $\int \frac{3x+7}{(x+2)^2(x+3)} dx$ .

First observe that the degree of the numerator is less than the degree of the denominator. In this case the correct form of the partial fraction expansion is

$$\frac{a}{(x + 2)} + \frac{b}{(x + 2)^2} + \frac{c}{(x + 3)}.$$

The reason there are two terms devoted to  $(x + 2)$  is that this is squared. Computing the constants yields

$$\frac{3x + 7}{(x + 2)^2(x + 3)} = \frac{1}{(x + 2)^2} + \frac{2}{x + 2} - \frac{2}{x + 3}$$

and therefore,

$$\int \frac{3x + 7}{(x + 2)^2(x + 3)} dx = -\frac{1}{x + 2} + 2 \ln(x + 2) - 2 \ln(x + 3) + C.$$

**Example 9.11.9** Find the proper form for the partial fractions expansion of

$$\frac{x^3 + 7x + 9}{(x^2 + 2x + 2)^3(x + 2)^2(x + 1)(x^2 + 1)}.$$

First check to see if the degree of the numerator is smaller than the degree of the denominator. Since this is the case, look for a partial fractions decomposition in the following form.

$$\frac{ax + b}{(x^2 + 2x + 2)} + \frac{cx + d}{(x^2 + 2x + 2)^2} + \frac{ex + f}{(x^2 + 2x + 2)^3} +$$

$$\frac{A}{(x + 2)} + \frac{B}{(x + 2)^2} + \frac{D}{(x + 1)} + \frac{gx + h}{x^2 + 1}.$$

These examples illustrate what to do when using the method of partial fractions. You first check to be sure the degree of the numerator is less than the degree of the denominator. If this is not so, do a long division. Then you factor the denominator into a product of factors, some linear of the form  $ax + b$  and others quadratic,  $ax^2 + bx + c$  which cannot be factored further. Next follow the procedure illustrated in the above examples.

## 9.12 Rational Functions Of Trig. Functions

**Example 9.12.1** Find  $\int \frac{\cos \theta}{1+\cos \theta} d\theta$ .

The integrand is an example of a rational function of cosines and sines. When such a thing occurs there is a substitution which will reduce the integrand to a rational function like those above which can then be integrated using partial fractions. The substitution is  $u = \tan\left(\frac{\theta}{2}\right)$ . Thus in this example,  $du = \left(1 + \tan^2\left(\frac{\theta}{2}\right)\right) \frac{1}{2} d\theta$  and so in terms of this new variable, the indefinite integral is

$$\int \frac{2 \cos(2 \arctan u)}{(1 + \cos(2 \arctan u))(1 + u^2)} du.$$

You can evaluate  $\cos(2 \arctan u)$  exactly. This equals  $2 \cos^2(\arctan u) - 1$ . Setting up a little triangle as above,  $\cos(\arctan u)$  equals  $1/\sqrt{1+u^2}$  and so the integrand reduces to

$$\frac{2 \left(2 \left(1/\sqrt{1+u^2}\right)^2 - 1\right)}{\left(1 + \left(2 \left(1/\sqrt{1+u^2}\right)^2 - 1\right)\right) (1 + u^2)} = \frac{1 - u^2}{1 + u^2} = -1 + \frac{2}{1 + u^2}$$

therefore, in terms of  $u$  the antiderivative equals  $-u + 2 \arctan u$ . Now replace  $u$  to obtain

$$-\tan\left(\frac{\theta}{2}\right) + 2 \arctan\left(\tan\left(\frac{\theta}{2}\right)\right) + C.$$

When you use partial fractions, **be sure you look for something which is of the right form**. Otherwise you may be looking for something which is not there.

## 9.13 Exercises

1. Give a condition on  $a$ ,  $b$ , and  $c$  such that  $ax^2 + bx + c$  cannot be factored as a product of two polynomials which have real coefficients.
2. Find the partial fractions expansion of the following rational functions.
  - (a)  $\frac{2x+7}{(x+1)^2(x+2)}$
  - (b)  $\frac{5x+1}{(x^2+1)(2x+3)}$
  - (c)  $\frac{5x+1}{(x^2+1)^2(2x+3)}$
  - (d)  $\frac{5x^4+10x^2+3+4x^3+6x}{(x+1)(x^2+1)^2}$
3. Find the antiderivatives
  - (a)  $\int \frac{x^5+4x^4+5x^3+2x^2+2x+7}{(x+1)^2(x+2)} dx$
  - (b)  $\int \frac{5x+1}{(x^2+1)(2x+3)} dx$
  - (c)  $\int \frac{5x+1}{(x^2+1)^2(2x+3)} dx$
4. Find  $\int \frac{\sin \theta}{1+\sin \theta} d\theta$ . **Hint:** Use the above procedure of letting  $u = \tan\left(\frac{\theta}{2}\right)$  and then multiply both the top and the bottom by  $(1 - \sin \theta)$  to see another way of doing it.
5. In finding  $\int \sec(x) dx$ , try the substitution  $u = \sin(x)$ .

6. In finding  $\int \csc(x) dx$  try the substitution  $u = \cos(x)$ .

7. Find the antiderivatives.

(a)  $\int \frac{17x-3}{(6x+1)(x-1)} dx$

(b)  $\int \frac{50x^4-95x^3-20x^2-3x+7}{(5x+3)(x-2)(2x-1)} dx$  **Hint:** Notice the degree of the numerator is larger than the degree of the denominator.

(c)  $\int \frac{8x^2+x-5}{(3x+1)(x-1)(2x-1)} dx$

(d)  $\int \frac{3x+2}{(5x+3)(x+1)} dx$

8. Find the antiderivatives

(a)  $\int \frac{52x^2+68x+46+15x^3}{(x+1)^2(5x^2+10x+8)} dx$

(b)  $\int \frac{9x^2-42x+38}{(3x+2)(3x^2-12x+14)} dx$

(c)  $\int \frac{9x^2-6x+19}{(3x+1)(3x^2-6x+5)} dx$

9. Find the antiderivatives.

(a)  $\int \frac{1}{(3x^2+12x+13)^2} dx$

(b)  $\int \frac{1}{(5x^2+10x+7)^2} dx$

(c)  $\int \frac{1}{(5x^2-20x+23)^2} dx$

## 9.14 Practice Problems For Antiderivatives

Here are lots of practice problems for finding antiderivatives. Some of these are very hard but you don't have to do all of them if you don't want to. However, the more you do, the better you will be at taking antiderivatives. Most of these problems are modifications of problems I found in a Russian calculus book. This book had some which were even harder. I shall give answers to these problems so you can see whether you have it right. Beware that sometimes you may get it right even though it looks different than the answer given. Also, there is no guarantee that my answers are right.

1. Find  $\int \frac{\sqrt{2x+1}}{x} dx$ .

Answer:

$$\int \frac{\sqrt{2x+1}}{x} dx = 2\sqrt{2}\sqrt{x} + \ln x + C$$

2. Find  $\int \frac{te^t}{(t+1)^2} dt$  **Hint:** Write this as

$$\int \left( \frac{(1+t)e^t}{(t+1)^2} - \frac{e^t}{(1+t)^2} \right) dt$$

$$= \int \left( \frac{e^t}{(t+1)} - \frac{e^t}{(1+t)^2} \right) dt.$$

Answer:  $\frac{e^t}{t+1} + C$

3. Find  $\int \frac{5-2x^2}{5+2x^2} dx$ .

Answer:

$$\int \frac{5-2x^2}{5+2x^2} dx = -x + \sqrt{10} \arctan \frac{1}{5}x\sqrt{10} + C$$

4. Find  $\int (3-x^5)^2 dx$ .

Answer:

$$\int (3-x^5)^2 dx = \frac{1}{11}x^{11} - x^6 + 9x + C$$

5. Find  $\int (2x+3)^{-25} dx$ .

Answer:

$$\int (2x+3)^{-25} dx = -\frac{1}{48(2x+3)^{24}} + C$$

6. Find  $\int 5^x dx$ .

Answer:

$$\int 5^x dx = \frac{1}{\ln 5} 5^x + C$$

7. Find  $\int \cosh^2 8x dx$ .

Answer:

$$\int \cosh^2 8x dx = \frac{1}{16} \cosh 8x \sinh 8x + \frac{1}{2}x + C$$

8. Find  $\int \tanh^2 2x dx$ .

Answer:

$$\int \tanh^2 2x dx =$$

$$-\frac{1}{2} \tanh 2x - \frac{1}{4} \ln(-1 + \tanh 2x)$$

$$+ \frac{1}{4} \ln(1 + \tanh 2x) + C$$

9. Find  $\int \cosh(3x+3) dx$ .

Answer:

$$\int \cosh(3x+3) dx = \frac{1}{3} \sinh(3x+3) + C$$

10. Find  $\int 8^{1+x} dx$ .

Answer:

$$\int 8^{1+x} dx = \frac{2}{\ln 2} 2^{3x} + \frac{2}{3\ln 2} 2^{3x} + C$$

11. Find  $\int \frac{\sqrt{(36+x^2)} + \sqrt{(36-x^2)}}{\sqrt{(1296-x^4)}} dx$ .

Answer:

$$\int \frac{\sqrt{(36+x^2)} + \sqrt{(36-x^2)}}{\sqrt{(1296-x^4)}} dx = \int \frac{1}{\sqrt{(36-x^2)}} dx +$$

$$\int \frac{1}{\sqrt{(36+x^2)}} dx = \arcsin \frac{1}{6}x +$$

$$\ln \left( x + \sqrt{(36+x^2)} \right) + C$$

12. Find  $\int (7x+1)^{30} dx$ .

Answer:

$$\int (7x+1)^{30} dx = \frac{1}{217} (7x+1)^{31} + C$$

13. Find  $\int \sqrt{(1+\sin 3x)} dx$ .

Answer:

$$\int \sqrt{(1+\sin 3x)} dx =$$

$$\frac{2}{3} (\sin 3x - 1) \frac{\sqrt{(1+\sin 3x)}}{\cos 3x} + C$$

14. Find  $\int \left( \sqrt{(3+2x)} \right)^5 dx$ .

Answer:

$$\int \left( \sqrt{(3+2x)} \right)^5 dx = \frac{1}{7} \left( \sqrt{(3+2x)} \right)^7 + C$$

15. Find  $\int \frac{1}{49+4x^2} dx$ .

Answer:

$$\int \frac{1}{49+4x^2} dx = \frac{1}{14} \arctan \frac{2}{7}x + C$$

16. Find  $\int \frac{1}{\sqrt{(4x^2-9)}} dx$ .

Answer:

$$\int \frac{1}{\sqrt{(4x^2-9)}} dx =$$

$$\frac{1}{2} \ln \left( 2x + \sqrt{(4x^2-9)} \right) + C$$

17. Find  $\int \frac{1}{\sin^2(2x+3)} dx$ .

Answer:

$$\int \frac{1}{\sin^2(2x+3)} dx =$$

$$-\frac{1}{2\sin(2x+3)} \cos(2x+3) + C$$

18. Find  $\int \frac{1}{1+\cos 6x} dx$ .

Answer:

$$\int \frac{1}{1+\cos 6x} dx = \left( \tan \frac{1}{2} \right) x + C$$

19. Find  $\int \frac{1}{1+\sin(3x)} dx$ .

Answer:

$$\int \frac{1}{1+\sin 3x} dx = -\frac{2}{3(\tan \frac{3}{2}x+1)} + C$$

20. Find  $\int x^2 \sqrt{(4x^3 + 2)} dx$ .

Answer:

$$\int x^2 \sqrt{(4x^3 + 2)} dx = \frac{1}{18} \left( \sqrt{(4x^3 + 2)} \right)^3 + C$$

21. Find  $\int x^8 \left( \sqrt{(3x^9 + 2)} \right)^5 dx$ .

Answer:

$$\begin{aligned} \int x^8 \left( \sqrt{(3x^9 + 2)} \right)^5 dx \\ = \frac{2}{189} \left( \sqrt{(3x^9 + 2)} \right)^7 + C \end{aligned}$$

22. Find  $\int \frac{x}{3+8x^4} dx$ .

Answer:

$$\begin{aligned} \int \frac{x}{3+8x^4} dx = \\ \frac{1}{24} \sqrt{6} \arctan \left( \frac{2}{3} x^2 \sqrt{6} \right) + C \end{aligned}$$

23. Find  $\int \frac{1}{3(2+x)\sqrt{x}} dx$ .

Answer:

$$\begin{aligned} \int \frac{1}{3(2+x)\sqrt{x}} dx = \\ \frac{1}{9} \sqrt{3} \sqrt{6} \arctan \left( \frac{1}{6} \sqrt{3} \sqrt{x} \sqrt{6} \right) + C \end{aligned}$$

24. Find  $\int \frac{1}{x\sqrt{(25x^2-9)}} dx$ .

Answer:

You could let  $3 \sec u = 5x$  so

$(3 \sec u \tan u) du = 5 dx$  and then

$$\int \frac{1}{x\sqrt{(25x^2-9)}} dx = \frac{1}{3} \int du = \frac{u}{3} + C.$$

Now restoring the original variables, this yields  $\frac{1}{3} \operatorname{arcsec} \frac{5}{3} |x| + C$ .

25. Find  $\int \frac{1}{\sqrt{(x(3x-5))}} dx$ .

Answer:

$$\begin{aligned} \int \frac{dx}{\sqrt{(x(3x-5))}} = \\ \frac{1}{3} \sqrt{3} \ln \left( \sqrt{3} \left( x - \frac{5}{6} \right) + \sqrt{(3x^2 - 5x)} \right) + C \end{aligned}$$

26. Find  $\int \frac{1}{x \cos(\ln 4x)} dx$ .

Answer:

$$\begin{aligned} \int \frac{1}{x \cos(\ln 4x)} dx = \\ \ln (\sec (\ln 4x) + \tan (\ln 4x)) + C \end{aligned}$$

27. Find  $\int x^7 e^{x^8} dx$ .

Answer:

$$\int x^7 e^{x^8} dx = \frac{1}{8} e^{x^8} + C$$

28. Find  $\int \frac{\ln^4 x}{x} dx$ .

Answer:

$$\int \frac{\ln^4 x}{x} dx = \frac{1}{5} \ln^5 x + C$$

29. Find  $\int \frac{1}{x(2+\ln 2x)} dx$ .

Answer:

$$\begin{aligned} \int \frac{1}{x(2+\ln 2x)} dx = \\ \ln (2 + \ln 2x) + C \end{aligned}$$

30. Find  $\int \frac{\sin 7x + \cos 7x}{\sqrt{(\sin 7x - \cos 7x)}} dx$ .

Answer:

$$\begin{aligned} \int \frac{\sin 7x + \cos 7x}{\sqrt{(\sin 7x - \cos 7x)}} dx = \\ \frac{2}{7} \sqrt{(\sin 7x - \cos 7x)} + C \end{aligned}$$

31. Find  $\int \csc 4x dx$ .

Answer:

$$\begin{aligned} \int \csc 4x dx = \\ \frac{1}{4} \ln |\csc 4x - \cot 4x| + C \end{aligned}$$

32. Find  $\int \frac{\arctan 5x}{1+25x^2} dx$ .

Answer:

$$\int \frac{\arctan 5x}{1+25x^2} dx = \frac{1}{10} \arctan^2 5x + C$$

33. Find  $\int \frac{18}{(6+2x)(6-x)} \cos \left( \ln \frac{6+2x}{6-x} \right) dx$ .

Answer:

$$\begin{aligned} \int \frac{18}{(6+2x)(6-x)} \cos \left( \ln \frac{6+2x}{6-x} \right) dx = \\ \sin \left( \ln \frac{6+2x}{6-x} \right) + C \end{aligned}$$

34. Find  $\int x^{23} (2 - 6x^{12})^{10} dx$

Answer:

$$\begin{aligned} \int x^{23} (2 - 6x^{12})^{10} dx = \\ \frac{1}{5184} (2 - 6x^{12})^{12} - \frac{1}{2376} (2 - 6x^{12})^{11} + C \end{aligned}$$

35. Find  $\int \frac{x^5}{\sqrt{(6-3x^2)}} dx$ .

Answer:

$$\begin{aligned} \int \frac{x^5}{\sqrt{(6-3x^2)}} dx = \\ -\frac{1}{15} x^4 \sqrt{(6-3x^2)} - \frac{8}{45} x^2 \sqrt{(6-3x^2)} \\ - \frac{32}{45} \sqrt{(6-3x^2)} + C \end{aligned}$$

36. Find  $\int \cos^3 3x \sin^{\frac{1}{2}} 3x \, dx$ .

Answer:

$$\begin{aligned} \int \cos^3 3x \sin^{\frac{1}{2}} 3x \, dx &= \\ \int \left( \cos 3x (1 - \sin^2 3x) \sin^{\frac{1}{2}} 3x \right) dx &= \\ = -\frac{2}{21} \sin^{\frac{7}{2}} 3x + \frac{2}{9} \sin^{\frac{3}{2}} 3x + C \end{aligned}$$

37. Find  $\int \frac{1}{e^{2x} + e^x} \, dx$ .

Answer:

$$\int \frac{1}{e^{2x} + e^x} \, dx = \frac{-1 - xe^x}{e^x} + \ln(e^x + 1) + C$$

38. Find  $\int \frac{1}{\sqrt{(e^x + 1)}} \, dx$ .

Answer:

Let  $u^2 = e^x$  so  $2u \, du = e^x \, dx = au^2 \, dx$ .  
In terms of  $u$  this is

$$2 \int \frac{1}{u\sqrt{1+u^2}} \, du. \text{ Now let}$$

$u = \tan \theta$  so  $du = (\sec^2 \theta) \, d\theta$ . Then the indefinite integral becomes

$$\begin{aligned} 2 \int \frac{\sec^2 \theta}{\tan \theta \sec(\theta)} \, d\theta &= 2 \int \csc \theta \, d\theta = \\ 2 \ln |\csc \theta - \cot \theta| + C. \end{aligned}$$

In terms of  $u$  this is

$$2 \ln \left| \frac{\sqrt{u^2 + 1}}{u} - \frac{1}{u} \right| + C \text{ and in terms of } x$$

this is

$$2 \ln \left| \frac{\sqrt{(e^x + 1)}}{e^{\frac{1}{2}x}} - e^{-\frac{1}{2}x} \right| + C.$$

39. Find  $\int \frac{\arctan \sqrt{x}}{\sqrt{x(1+x)}} \, dx$ .

Answer:

$$\int \frac{\arctan \sqrt{x}}{\sqrt{x(1+x)}} = \arctan^2 \sqrt{x} + C$$

40. Find  $\int \sqrt{\left(\frac{1+x}{1-x}\right)} \, dx$ .

Answer:

Multiply the fraction on the top and bottom by  $1 + x$  to get

$$\int \frac{1+x}{\sqrt{(1-x^2)}} \, dx. \text{ Now let } x = \sin \theta \text{ so } dx = \cos \theta \, d\theta. \text{ Then this is}$$

$$\begin{aligned} \int \frac{1+\sin \theta}{\sqrt{(1-\sin^2 \theta)}} \cos \theta \, d\theta &= \\ \int (1 + \sin \theta) \, d\theta &= \theta - \cos \theta + C. \end{aligned}$$

In terms of  $x$  this gives

$$\arcsin x - \sqrt{(1-x^2)} + C.$$

41. Find  $\int \sqrt{\left(\frac{x-2}{x+2}\right)} \, dx$ .

Answer:

Multiply the fraction on the top and bottom by  $x + 2$  to get  $\int \frac{\sqrt{(x^2-4)}}{x+2} \, dx$   
Now let  $x = 2 \sec \theta$  so  $dx = 2 \sec \theta \tan \theta \, d\theta$   
and in terms of  $\theta$  the indefinite integral is

$$\begin{aligned} 2 \int \frac{\sqrt{\sec^2(\theta)-1}}{\sec \theta + 1} \sec \theta \tan \theta \, d\theta &= \\ 2 \int \frac{\tan^2 \theta \sec(\theta)}{1 + \sec \theta} \, d\theta &= \\ = 2 \int \frac{\tan^2 \theta}{1 + \cos \theta} = 2 \int \sec^2 \theta \, d\theta - 2 \int \sec \theta \, d\theta &= \\ 2 \tan \theta - 2 \ln(\sec \theta + \tan \theta) + C. \text{ Now in} & \\ \text{terms of } x \text{ this is} & \end{aligned}$$

$$\begin{aligned} \sqrt{(x^2-4)} - \\ 2 \ln \left( \frac{1}{2}x + \frac{1}{2}\sqrt{(x^2-4)} \right) + C. \end{aligned}$$

42. Find  $\int \frac{x^2}{\sqrt{(4+9x^2)}} \, dx$ .

Answer:

$$\begin{aligned} \int \frac{x^2}{\sqrt{(4+9x^2)}} \, dx &= \\ \frac{1}{18}x\sqrt{(4+9x^2)} - \\ \frac{2}{27} \ln \left( 3x + \sqrt{(4+9x^2)} \right) + C \end{aligned}$$

43. Find  $\int \frac{1}{\sqrt{(x^2+49)}} \, dx$ .

Answer:

$$\int \frac{1}{\sqrt{(x^2+49)}} \, dx = \ln \left( x + \sqrt{(x^2+49)} \right) + C$$

44. Find  $\int \frac{1}{\sqrt{(x^2-36)}} \, dx$ .

Answer:

$$\int \frac{1}{\sqrt{(x^2-36)}} \, dx = \ln \left( x + \sqrt{(x^2-36)} \right) + C$$

45. Find  $\int x \ln(3x) \, dx$ .

Answer:

$$\int x \ln(3x) \, dx = \frac{1}{2}x^2 \ln 3x - \frac{1}{4}x^2 + C$$

46. Find  $\int x \ln^2(6x) \, dx$ .

Answer:

$$\int x \ln^2 6x \, dx = \frac{1}{2}x^2 \ln^2 6x - \frac{1}{2}x^2 \ln 6x + \frac{1}{4}x^2 + C$$

47. Find  $\int x^3 e^{x^2} dx$ .

Answer:

$$\int x^3 e^{x^2} dx = \frac{1}{2} x^2 e^{x^2} - \frac{1}{2} e^{x^2} + C$$

48. Find  $\int x^2 \sin 8x dx$

Answer:

$$\begin{aligned} \int x^2 \sin 8x dx &= \\ -\frac{1}{8} x^2 \cos 8x + \frac{1}{256} \cos 8x + \frac{1}{32} x \sin 8x + C \end{aligned}$$

49. Find  $\int \arcsin x dx$

Answer:

$$\begin{aligned} \int \arcsin x dx &= \\ x \arcsin x + \sqrt{1-x^2} + C \end{aligned}$$

50. Find  $\int \arctan x dx$

Answer:

$$\begin{aligned} \int \arctan x dx &= \\ x \arctan x - \frac{1}{2} \ln(1+x^2) + C \end{aligned}$$

51. Find  $\int \frac{\arctan 6x}{x^3} dx$

Answer:

$$\begin{aligned} \int \frac{\arctan 6x}{x^3} dx &= \\ -\frac{1}{2x^2} \arctan 6x - \frac{3}{x} - 18 \arctan 6x + C \end{aligned}$$

52. Find  $\int \sin 4x \ln(\tan 4x) dx$

Answer:

$$\begin{aligned} \text{Integration by parts gives} \\ -\frac{1}{4} \cos 4x \ln(\tan 4x) + \\ \frac{1}{4} \ln(\csc 4x - \cot 4x) + C \end{aligned}$$

53. Find  $\int e^{2x} \sqrt{e^{4x} + 1} dx$ .

Answer:

$$\begin{aligned} \int e^{2x} \sqrt{e^{4x} + 1} dx &= \\ \frac{1}{4} e^{2x} \sqrt{e^{4x} + 1} + \frac{1}{4} \operatorname{arcsinh}(e^{2x}) + C \end{aligned}$$

54. Find  $\int \cos(\ln 7x) dx$ .

Answer:

$$\begin{aligned} \int \cos(\ln 7x) dx &= x \cos(\ln 7x) \\ + \int x \sin(\ln(7x)) \left(\frac{1}{x}\right) dx &= x \cos(\ln 7x) + \\ [x \sin(\ln 7x) - \int \cos(\ln(7x)) dx] \text{ and so} \\ \int \cos(\ln 7x) dx &= \\ \frac{1}{2} [x \cos(\ln 7x) + x \sin(\ln 7x)] + C \end{aligned}$$

55. Find  $\int \sin(\ln 7x) dx$ .

Answer:

$$\begin{aligned} \int \sin(\ln 7x) dx &= \\ \frac{1}{2} [x \sin(\ln 7x) - x \cos(\ln 7x)] + C \end{aligned}$$

56. Find  $\int e^{5x} \sin 3x dx$ .

Answer:

$$\begin{aligned} \int e^{5x} \sin 3x dx &= \\ -\frac{3}{34} e^{5x} \cos 3x + \frac{5}{34} e^{5x} \sin 3x + C \end{aligned}$$

57. Find  $\int e^{3x} \sin^2 2x dx$ .

Answer:

$$\begin{aligned} \int e^{3x} \sin^2 2x dx &= \\ \frac{1}{6} e^{3x} - \frac{3}{25} e^{3x} (\cos 2) x - \frac{4}{25} e^{3x} (\sin 2) x + C \end{aligned}$$

58. Find  $\int 3 \frac{x+4}{(x+3)^2} dx$ .

Answer:

$$\begin{aligned} \int 3 \frac{x+4}{(x+3)^2} dx &= \\ 3 \ln(x+3) - \frac{3}{x+3} + C \end{aligned}$$

59. Find  $\int \frac{7x^2+100x+347}{(x+2)(x+7)^2} dx$ .

Answer:

$$\begin{aligned} \int \frac{7x^2+100x+347}{(x+2)(x+7)^2} dx &= \\ 7 \ln(x+2) - \frac{2}{x+7} + C \end{aligned}$$

60. Find  $\int 2 \frac{3x^2+7x+8}{(x+2)(x^2+2x+2)} dx$ .

Answer:

$$\begin{aligned} \int 2 \frac{3x^2+7x+8}{(x+2)(x^2+2x+2)} dx &= \\ 6 \ln(x+2) + 2 \arctan(1+x) + C \end{aligned}$$

61. Find  $\int 2 \frac{4x^2+65x+266}{(x+2)(x^2+16x+66)} dx$ .

Answer:

$$\begin{aligned} \int 2 \frac{4x^2+65x+266}{(x+2)(x^2+16x+66)} dx &= \\ 8 \ln(x+2) + \\ \sqrt{2} \arctan \frac{1}{4} (16+2x) \sqrt{2} + C \end{aligned}$$

62. Find  $\int \frac{1}{x^3+8} dx$ .

Answer:

$$\begin{aligned} \int \frac{1}{x^3+8} dx &= \\ \frac{1}{12} \ln(x+2) - \frac{1}{24} \ln(x^2-2x+4) \\ + \frac{1}{12} \sqrt{3} \arctan \frac{1}{6} (2x-2) \sqrt{3} + C \end{aligned}$$

63. Find  $\int \frac{1}{x^4+x^2+1} dx$ .

Answer:

$$\begin{aligned} &\int \frac{1}{x^4+x^2+1} dx = \\ &\frac{1}{4} \ln(x^2 + x + 1) + \\ &\frac{1}{6} \sqrt{3} \arctan \frac{1}{3} (2x + 1) \sqrt{3} \\ &- \frac{1}{4} \ln(x^2 - x + 1) + \\ &\frac{1}{6} \sqrt{3} \arctan \frac{1}{3} (2x - 1) \sqrt{3} + C \end{aligned}$$

64. Find  $\int \frac{1}{x^4+81} dx$ . **Hint:**

$$\begin{aligned} x^4 + 81 &= \\ (x^2 - 3x\sqrt{2} + 9)(x^2 + 3x\sqrt{2} + 9). \end{aligned}$$

Answer:

Now you can use partial fractions to find

$$\begin{aligned} &\int \frac{1}{x^4+81} dx = \\ &\frac{1}{216} \sqrt{2} \ln \frac{x^2+3x\sqrt{2}+3}{x^2-3x\sqrt{2}+3} + \\ &\frac{1}{108} \sqrt{2} \arctan \left( \frac{1}{3} x \sqrt{2} + 1 \right) \\ &+ \frac{1}{108} \sqrt{2} \arctan \left( \frac{1}{3} x \sqrt{2} - 1 \right) + C \end{aligned}$$

65. Find  $\int \frac{1}{x^6+64} dx$ .

Answer:

First factor the denominator.

$$\begin{aligned} x^6 + 64 &= \\ (x^2 + 4) \left( (x - \sqrt{3})^2 + 1 \right) \left( (x + \sqrt{3})^2 + 1 \right). \end{aligned}$$

$\frac{1}{x^6+64} = \frac{Ax+B}{x^2+4} + \frac{Cx+D}{(x-\sqrt{3})^2+1} + \frac{Ex+F}{(x+\sqrt{3})^2+1}$  and so after wading through much affliction, the partial fractions decomposition is

$\frac{16}{3x^2+12} + \frac{-\frac{1}{192}\sqrt{3}x+\frac{1}{48}}{(x-\sqrt{3})^2+1} + \frac{\frac{1}{192}\sqrt{3}x+\frac{1}{48}}{(x+\sqrt{3})^2+1}$ . Therefore, the indefinite integral is

$$\begin{aligned} &\int \left( \frac{16}{3x^2+12} + \frac{-\frac{1}{192}\sqrt{3}x+\frac{1}{48}}{(x-\sqrt{3})^2+1} + \right. \\ &\left. + \frac{\frac{1}{192}\sqrt{3}x+\frac{1}{48}}{(x+\sqrt{3})^2+1} \right) dx = \\ &\frac{1}{96} \arctan \frac{1}{2} x - \frac{1}{384} \sqrt{3} \ln(x^2 - 2\sqrt{3}x + 4) \\ &+ \frac{1}{192} \arctan(x - \sqrt{3}) + \\ &\frac{1}{384} \sqrt{3} \ln(x^2 + 2\sqrt{3}x + 4) + \\ &\frac{1}{192} \arctan(x + \sqrt{3}) + C \end{aligned}$$

66. Find  $\int \frac{x^2}{(x-3)^{100}} dx$ .

Answer:

$$\begin{aligned} &\int \frac{x^2}{(x-3)^{100}} dx \\ &= \frac{1}{11(3-x)^{99}} - \frac{3}{49(3-x)^{98}} + \frac{1}{97(3-x)^{97}} + C \end{aligned}$$

67. Find  $\int 2 \frac{4+4x+3x^2+x^3}{(x^2+1)(x^2+2x+5)} dx$ .

Answer:

$$\begin{aligned} &\int 2 \frac{4+4x+3x^2+x^3}{(x^2+1)(x^2+2x+5)} dx \\ &= \frac{1}{2} \ln(x^2 + 1) + \arctan x + \\ &\frac{1}{2} \ln(x^2 + 2x + 5) + \arctan \left( \frac{1}{2} + \frac{1}{2}x \right) + C \end{aligned}$$

68. Find  $\int \frac{1-x^7}{x(1+x^7)} dx$

Answer:

$$\begin{aligned} &\int \frac{1-x^7}{x(1+x^7)} dx \\ &= \ln x - \frac{2}{7} \ln(1+x^7) + C \end{aligned}$$

69. Find  $\int \frac{x^2+1}{x^4-2x^2+1} dx$

Answer:

$$\begin{aligned} &\int \frac{x^2+1}{x^4-2x^2+1} dx \\ &= -\frac{1}{2(x-1)} - \frac{1}{2(1+x)} + C \end{aligned}$$

70. Find  $\int \frac{x^5}{x^8+1} dx$

Answer:

$$\begin{aligned} &\text{Let } u = x^2. \int \frac{x^5}{x^8+1} dx \\ &= \int \frac{u^2}{2u^4+2} du. \text{ Now use partial fractions.} \\ &\int \frac{x^5}{x^8+1} dx = \frac{1}{16} \sqrt{2} \ln \frac{x^4-x^2\sqrt{2}+1}{x^4+x^2\sqrt{2}+1} \\ &+ \frac{1}{8} \sqrt{2} \arctan(x^2\sqrt{2} + 1) \\ &+ \frac{1}{8} \sqrt{2} \arctan(x^2\sqrt{2} - 1) + C \end{aligned}$$

71. Find  $\int \frac{x^2+4}{x^6+64} dx$ .

Answer:

$$\begin{aligned} &\int \frac{x^2+4}{x^6+64} dx = \\ &-\frac{1}{96} \sqrt{3} \ln(x^2 - 2\sqrt{3}x + 4) + \\ &\frac{1}{16} \arctan(x - \sqrt{3}) + \\ &\frac{1}{96} \sqrt{3} \ln(x^2 + 2\sqrt{3}x + 4) \\ &+ \frac{1}{16} \arctan(x + \sqrt{3}) + C \end{aligned}$$



72. Find  $\int \frac{dx}{x(2+3\sqrt{x}+\sqrt[3]{x})}$ .

Answer:

$$\begin{aligned}\int \frac{dx}{x(2+3\sqrt{x}+\sqrt[3]{x})} &= \int \frac{6}{u(2+3u^3+u^2)} du = \\ &= 3 \ln u - \frac{6}{7} \ln(1+u) - \frac{15}{14} \ln(3u^2-2u+2) \\ &\quad - \frac{3}{35} \sqrt{5} \arctan \frac{1}{10} (6u-2) \sqrt{5} + C \\ &= 3 \ln x^{1/6} - \frac{6}{7} \ln(1+x^{1/6}) - \\ &\quad \frac{15}{14} \ln(3x^{1/3}-2x^{1/6}+2) - \\ &\quad \frac{3}{35} \sqrt{5} \arctan \frac{1}{10} (6x^{1/6}-2) \sqrt{5} + C\end{aligned}$$

73. Find  $\int \frac{\sqrt{(x+3)}-\sqrt{(x-3)}}{\sqrt{(x+3)}+\sqrt{(x-3)}} dx$ .

Answer:

$$\begin{aligned}\int \frac{\sqrt{(x+3)}-\sqrt{(x-3)}}{\sqrt{(x+3)}+\sqrt{(x-3)}} dx &= \\ &= \frac{1}{6} x^2 - \frac{1}{6} \sqrt{(x+3)} \left( \sqrt{(x-3)} \right)^3 \\ &\quad - \frac{1}{2} \sqrt{(x+3)} \sqrt{(x-3)} + \\ &\quad \frac{3}{2} \frac{\sqrt{((x-3)(x+3))}}{\sqrt{(x-3)}\sqrt{(x+3)}} \ln \left( x + \sqrt{(x^2-9)} \right) + C\end{aligned}$$

74. Find  $\int \frac{x^2}{\sqrt{x^2+6x+13}} dx$ .

Answer:

$$\begin{aligned}\int \frac{x^2}{\sqrt{x^2+6x+13}} dx &= \\ &= \frac{1}{2} x \sqrt{(x^2+6x+13)} - \frac{9}{2} \sqrt{(x^2+6x+13)} + \\ &\quad 7 \ln \left( x+3+\sqrt{(x^2+6x+13)} \right) + C\end{aligned}$$

75. Find  $\int \frac{\sqrt{(x^2+4x+5)}}{x} dx$ .

Answer:

$$\begin{aligned}\int \frac{\sqrt{(x^2+4x+5)}}{x} dx &= \\ &= \sqrt{(x^2+4x+5)} + 2 \operatorname{arcsinh}(x+2) \\ &\quad - \sqrt{5} \operatorname{arctanh} \frac{1}{5} (5+2x) \frac{\sqrt{5}}{\sqrt{(x^2+4x+5)}} + C\end{aligned}$$

You might try letting  $u = \sinh^{-1}(x+2)$ .

76. Find  $\int \frac{1}{x^3 \sqrt{x^2+9}} dx$ .

Answer:

$$\begin{aligned}\int \frac{1}{x^3 \sqrt{x^2+9}} dx &= \\ &= -\frac{1}{18x^2} \sqrt{(x^2+9)} + \\ &\quad \frac{1}{54} \operatorname{arctanh} \frac{3}{\sqrt{(x^2+9)}} + C\end{aligned}$$

You might try letting  $u = \frac{1}{\sqrt{x^2+9}}$ .

77. Find  $\int \frac{dx}{x^4 \sqrt{x^2-9}}$ .

Answer:

$$\begin{aligned}\int \frac{dx}{x^4 \sqrt{x^2-9}} &= \\ &= \frac{1}{27x^3} \sqrt{(x^2-9)} + \frac{2}{243x} \sqrt{(x^2-9)} + C\end{aligned}$$

78. Find  $\int \sin^6(3x) dx$ .

Answer:

$$\begin{aligned}\int \sin^6(3x) dx &= \\ &= -\frac{1}{2} \sin^5 x \cos x - \frac{5}{8} \sin^3 x \cos x \\ &\quad - \frac{15}{16} \cos x \sin x + \frac{5}{16} x + C\end{aligned}$$

79. Find  $\int \frac{\sin^3(x)}{\cos^4(x)} dx$ .

Answer:

$$\begin{aligned}\int \frac{\sin^3(x)}{\cos^4(x)} dx &= \\ &= \frac{1}{3} \frac{\sin^4 x}{\cos^3 x} - \frac{1}{3} \frac{\sin^4 x}{\cos x} - \\ &\quad \frac{1}{3} \sin^2 x \cos x - \frac{2}{3} \cos x + C\end{aligned}$$

80. Find  $\int \tan^5(2x) dx$ .

Answer:

$$\begin{aligned}\int \tan^5(2x) dx &= \\ &= -\frac{1}{4} \tan^2 2x + \frac{1}{8} \tan^4 2x \\ &\quad + \frac{1}{4} \ln(2+2\tan^2 2x) + C\end{aligned}$$

81. Find  $\int \frac{dx}{\sqrt{\tan(2x)}}$ .

Answer:

$$\begin{aligned}\int \frac{dx}{\sqrt{\tan(2x)}} &= \\ &= 2 \frac{\tan^{\frac{1}{2}} x}{1+2 \tan^2 x} + 2 \frac{\tan^{\frac{5}{2}} x}{1+2 \tan^2 x} \\ &\quad - 2 \tan^{\frac{1}{2}} x + \frac{1}{4} \sqrt{2} \arctan 2\sqrt{2} \frac{\tan^{\frac{1}{2}} x}{1-2 \tan x} \\ &\quad + \frac{1}{4} \sqrt{2} \ln \frac{2 \tan x + 2\sqrt{2} \tan^{\frac{1}{2}} x + 1}{\sqrt{(1+2 \tan^2 x)}} + C\end{aligned}$$

You might try the substitution,  $u = \tan(x)$ .

82. Find  $\int \frac{dx}{\cos x + 3 \sin x + 4}$ .

Answer:

$$\begin{aligned}\int \frac{dx}{\cos x + 3 \sin x + 4} &= \\ &= \frac{1}{3} \sqrt{6} \arctan \frac{1}{12} (6 \tan \frac{1}{2} x + 6) \sqrt{6} + C\end{aligned}$$

Try the substitution,  $u = \tan(\frac{x}{2})$ .

## 9.15 Volumes

Imagine a building having height  $h$  and consider the intersection of the building with a plane which is parallel to the ground at height  $y$ . Let the area of this intersection, called the cross section at height  $y$ , equal  $A(y)$ . Then the volume of the building between the two parallel planes at height  $y$  and height  $y + \Delta y$  would be approximately  $\Delta y (A(y))$ . Written in terms of differentials  $dV = A(y) dy$  and so the total volume of the building between 0 and  $y$ ,  $V(y)$ , would satisfy the differential equation,

$$\frac{dV}{dy} = A(y), \quad V(0) = 0.$$

The volume of the building would be  $V(h)$ . You can use differential equations to compute volumes of other shapes as well.

**Example 9.15.1** Consider a pyramid which sits on a square base of length 500 feet and suppose the pyramid has height 300 feet. Find the volume of the pyramid in cubic feet.

At height,  $y$ , the length of one of the sides,  $x$ , would satisfy  $\frac{x}{300-y} = \frac{500}{300} = \frac{5}{3}$  and so  $x = \frac{5}{3}(300 - y)$ . Therefore,

$$\frac{dV}{dy} = \left( \frac{5}{3}(300 - y) \right)^2$$

and so  $V(y) = -\frac{1}{5}(500 - \frac{5}{3}y)^3 + C$ . Since  $V(0) = 0$ , the constant must equal  $\frac{1}{5}(500)^3$ . Therefore, the total volume of the pyramid would equal

$$-\frac{1}{5} \left( 500 - \frac{5}{3}300 \right)^3 + \frac{1}{5}(500)^3 = 25,000,000 \text{ cubic feet.}$$

**Example 9.15.2** The base of a solid is a circle of radius 10 meters in the  $xy$  plane. When this solid is cut with a plane which is perpendicular to the  $xy$  plane and the  $x$  axis at  $x$ , the result is a square. Find the volume of the resulting solid.

Reasoning as above for the building and letting  $V(x)$  denote the volume between  $-10$  and  $x$ , yields

$$\frac{dV}{dx} = A(x), \quad V(-10) = 0$$

as an appropriate differential equation for this volume. Here  $A(x)$  is the area of the surface resulting from the intersection of the plane with the solid. The length of one side of this surface is  $2\sqrt{100 - x^2}$  because the equation of the circle which bounds the base is  $x^2 + y^2 = 100$ . Therefore,  $A(x) = 4(100 - x^2)$  and so  $V(x) = 400x - \frac{4}{3}x^3 + C$ . To find  $C$ , use  $V(-10) = 0$  and so  $C = -400(-10) + \frac{4}{3}(-10)^3 = \frac{8000}{3}$ . Therefore, the total volume is

$$V(10) = 400(10) - \frac{4}{3}(10)^3 + \frac{8000}{3} = \frac{16000}{3}$$

**Example 9.15.3** Find the volume of a sphere of radius  $R$ .

The sphere is obtained by revolving a disk of radius  $R$  about the  $y$  axis. Thus the radius of a the cross section at height  $y$  would be  $\sqrt{R^2 - y^2}$  and so the area of this cross section is  $\pi(R^2 - y^2)$ . Therefore,

$$\frac{dV}{dy} = \pi(R^2 - y^2), \quad V(-R) = 0$$

and so  $V(y) = \pi \left( R^2 y - \frac{1}{3} y^3 \right) + C$ . Now since  $V(-R) = 0$ ,  $C = \pi \left( -\frac{R^3}{3} + R^3 \right)$  and so

$$V(R) = \pi \left( R^3 - \frac{1}{3} (R)^3 \right) + \pi \left( -\frac{R^3}{3} + R^3 \right) = \frac{4}{3} \pi R^3$$

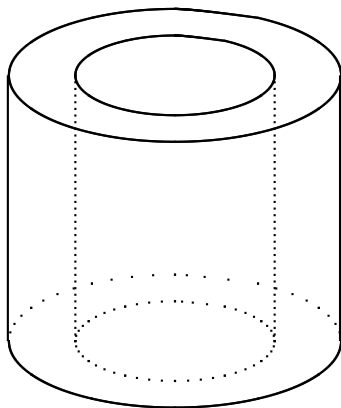
**Example 9.15.4** The graph of the function,  $f(x) = \sqrt{x}$  is revolved about the  $x$  axis. Find the volume of the resulting shape between  $x = 0$  and  $x = 8$ .

The cross sections perpendicular to the  $x$  axis for this shape are circles and the cross section at  $x$  has radius  $\sqrt{x}$ . Therefore,  $A(x) = \pi x$  and the appropriate differential equation and initial value is

$$\frac{dV}{dx} = \pi x, \quad V(0) = 0$$

and the answer is  $V(8)$ . From the differential equation and initial condition,  $V(x) = \pi \frac{x^2}{2}$  and so  $V(8) = \pi \frac{64}{2} = 32\pi$ .

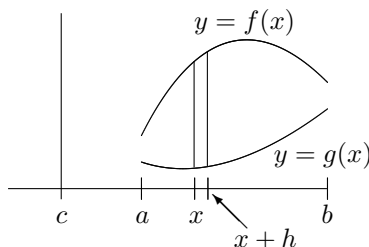
There is another way to find volumes without using cross sections. This method involves the notion of shells. Consider the following picture of a circular shell.



In this picture the radius of the inner circle will be  $r$  and the radius of the outer circle will be  $r + \Delta r$  while the height of the shell is  $H$ . Therefore, the volume of the shell would be the difference in the volumes of the two cylinders or

$$\pi(r + \Delta r)^2 H - \pi r^2 h = 2\pi H r (\Delta r) + \pi H (\Delta r)^2. \quad (9.12)$$

Now consider the problem of revolving the region between  $y = f(x)$  and  $y = g(x)$  for  $x \in [a, b]$  about the line  $x = c$  for  $c < a$ . The following picture is descriptive of the situation.



Let  $V(x)$  denote the volume of the solid which results from revolving the region between the graphs of  $f$  and  $g$  above the interval,  $[a, x]$  about the line  $x = c$ . Thus  $V(x+h) - V(x)$  equals the volume which results from revolving the region between the graphs of  $f$  and  $g$  which is also between the two vertical lines shown in the above picture. This results in a solid which is very nearly a circular shell like the one shown in the previous picture and the approximation gets better as let  $h$  decreases to zero. Therefore,

$$\begin{aligned} V'(x) &= \lim_{h \rightarrow 0} \frac{V(x+h) - V(x)}{h} \\ &= \lim_{h \rightarrow 0} \frac{2\pi |f(x) - g(x)| (x-c)h + \pi |f(x) - g(x)| h^2}{h} \\ &= 2\pi |f(x) - g(x)| (x-c). \end{aligned} \quad (9.13)$$

Also,  $V(a) = 0$  and this means that to find the volume of revolution it suffices to solve the initial value problem,

$$\frac{dV}{dx} = 2\pi |f(x) - g(x)| (x-c), \quad V(a) = 0 \quad (9.14)$$

and the volume of revolution will equal  $V(b)$ . Note that in the above formula, it is not necessary to worry about which is larger,  $f(x)$  or  $g(x)$  because it is expressed in terms of the absolute value of their difference. However, in doing the computations necessary to solve a given problem, you typically will have to worry about which is larger.

**Example 9.15.5** Find the volume of the solid formed by revolving the region between  $y = \sin(x)$  and the  $x$  axis for  $x \in [0, \pi]$  about the  $y$  axis.

In this example,  $c = 0$  and since  $\sin(x) \geq 0$  for  $x \in [0, \pi]$ , the initial value problem is

$$\frac{dV}{dx} = 2\pi x \sin(x), \quad V(0) = 0.$$

Then using integration by parts,

$$V(x) = 2\pi (\sin x - x \cos x) + C$$

and from the initial condition,  $C = 0$ . Therefore, the volume is  $V(\pi) = 2\pi^2$ .

**Example 9.15.6** Find the volume of the solid formed by revolving the region between  $y = \sin x$ ,  $y = \cos x$  and the  $x$  axis for  $x \in [0, \pi/4]$  about the line  $x = -4$

In this example,  $\cos x > \sin x$  for  $x \in [0, \pi/4]$  and so the initial value problem is

$$\frac{dV}{dx} = 2\pi (x+4) (\cos x - \sin x), \quad V(0) = 0.$$

Using integration by parts

$$V \in \int 2\pi (x+4) (\cos x - \sin x) dx = 2(5 \cos x + x \sin x + 3 \sin x + x \cos x) \pi + C$$

and the initial condition gives  $C = -10\pi$ . Therefore, the volume is

$$\begin{aligned} V(\pi/4) &= 2(5 \cos(\pi/4) + (\pi/4) \sin(\pi/4) + 3 \sin(\pi/4) + (\pi/4) \cos(\pi/4)) \pi - 10\pi \\ &= 2 \left( 4\sqrt{2} + \frac{1}{4}\pi\sqrt{2} \right) \pi - 10\pi \end{aligned}$$

**Example 9.15.7** Find the volume of a sphere of radius  $R$  using the method of shells.

In this case the volume of the sphere is obtained by revolving the region between  $y = \sqrt{R^2 - x^2}$ ,  $y = -\sqrt{R^2 - x^2}$  and the  $x$  axis for  $x \in [0, R]$  about the  $y$  axis. Thus this volume satisfies the initial value problem

$$\frac{dV}{dx} = 2\pi \left( 2\sqrt{R^2 - x^2} \right) x, \quad V(0) = 0.$$

Thus, using the method of substitution,

$$\int 2\pi \left( 2\sqrt{R^2 - x^2} \right) x \, dx = -\frac{4}{3} \left( \sqrt{R^2 - x^2} \right)^3 \pi + C$$

and from the initial condition,  $C = \frac{4}{3}R^3\pi$ . Therefore,

$$V(x) = -\frac{4}{3} \left( \sqrt{R^2 - x^2} \right)^3 \pi + \frac{4}{3}R^3\pi$$

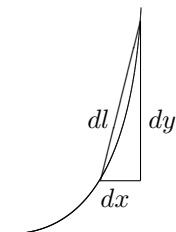
and so the volume of the sphere is  $V(R) = \frac{4}{3}R^3\pi$ , the same formula obtained earlier using the method of cross sections to set up the differential equation.

## 9.16 Exercises

1. The equation of an ellipse is  $\left(\frac{x}{a}\right)^2 + \left(\frac{y}{b}\right)^2 = 1$ . Sketch the graph of this in the case where  $a = 2$  and  $b = 3$ . You should get an oval shaped thing which looks like a squashed circle. Now find the volume of the solid obtained by revolving this ellipse about the  $y$  axis. What is the volume if it is revolved about the  $x$  axis? What is the volume of the solid obtained by revolving it about the line  $x = -2$ ?
2. A sphere of radius  $R$  has a hole drilled through a diameter which is centered at the diameter and of radius  $r < R$ . Find the volume of what is left after the hole has been drilled. How much material was taken out?
3. Show the volume of a right circular cone is  $(1/3) \times \text{area of the base} \times \text{height}$ .
4. Let  $R$  be a region in the  $xy$  plane of area  $A$  and consider the cone formed by fixing a point in space  $h$  units above  $R$  and taking the union of all lines starting at this point which end in  $R$ . Show that under these general conditions, the volume of the cone is  $(1/3) \times A \times h$ . **Hint:** The cross sections at height  $y$  look just like  $R$  but shrunk. Argue that the area at height  $y$ , denoted by  $A(y)$  is simply  $A(y) = A \frac{(h-y)^2}{h^2}$ .
5. A circle of radius  $r$  in the  $xy$  plane is the base of a solid which has the property that cross sections perpendicular to the  $x$  axis are equilateral triangles. Find the volume of this solid.
6. A square having each side equal to  $r$  in the  $xy$  plane is the base of a solid which has the property that cross sections perpendicular to the  $x$  axis are equilateral triangles. Find the volume of this solid.
7. The bounded region between  $y = x^2$  and  $y = x$  is revolved about the  $x$  axis. Find the volume of the solid which results.
8. The region between  $y = \ln x$ , and the  $x$  axis for  $x \in [1, 3]$  is revolved about the  $y$  axis. What is the volume of the resulting solid?
9. The region between  $y = 2 + \sin 3x$  and the  $x$  axis for  $x \in [0, \pi/3]$  is revolved about the line  $x = -1$ . Find the volume of the solid which results.

## 9.17 Lengths And Areas Of Surfaces Of Revolution

The same techniques can be used to compute lengths of the graph of a function,  $y = f(x)$ . Consider the following picture.



which depicts a small right triangle attached as shown to the graph of a function,  $y = f(x)$  for  $x \in [a, b]$ . If the triangle is small enough, this shows the length of the curve joined by the hypotenuse of the right triangle is essentially equal to the length of the hypotenuse. Thus,  $(dl)^2 = (dx)^2 + (dy)^2$  and dividing by  $(dx)^2$  yields

$$\frac{dl}{dx} = \sqrt{1 + \left(\frac{dy}{dx}\right)^2} = \sqrt{1 + f'(x)^2}, l(a) = 0$$

as an initial value problem for the function,  $l(x)$  which gives the length of this curve on  $[a, x]$ .

This definition gives the right answer for the length of a straight line. To see this, consider a straight line through the points  $(a, b)$  and  $(c, d)$  where  $a < c$ . Then the right answer is given by the Pythagorean theorem or distance formula and is  $\sqrt{(a - c)^2 + (b - c)^2}$ . What is obtained from the above initial value problem? The equation of the line is  $f(x) = b + \left(\frac{d-b}{c-a}\right)(x - a)$  and so  $f'(x) = \left(\frac{d-b}{c-a}\right)$ . Therefore, letting  $l$  denote the arc length function,

$$\frac{dl}{dx} = \sqrt{1 + \left(\frac{d-b}{c-a}\right)^2}, l(a) = 0.$$

Thus  $l(x) = \sqrt{1 + \left(\frac{d-b}{c-a}\right)^2}(x - a)$  and in particular,  $l(c)$ , the length of the line is given by

$$\sqrt{1 + \left(\frac{d-b}{c-a}\right)^2}(c - a) = \sqrt{(a - c)^2 + (b - c)^2}$$

as hoped. Thus this differential equation gives the right answer in the familiar cases but it also can be used to find lengths for more general curves than straight lines. Here is another familiar example.

**Example 9.17.1** Find the length of the part of the circle having radius  $r$  which is between the points  $\left(0, \frac{\sqrt{2}}{2}r\right)$  and  $\left(\frac{\sqrt{2}}{2}r, \frac{\sqrt{2}}{2}r\right)$ .

Here the function is  $f(x) = \sqrt{r^2 - x^2}$  and so  $f'(x) = -x/\sqrt{r^2 - x^2}$ . Therefore, our differential equation is

$$\frac{dl}{dx} = \sqrt{1 + \frac{x^2}{(r^2 - x^2)}} = \sqrt{\frac{r^2}{r^2 - x^2}} = \frac{r}{\sqrt{r^2 - x^2}}.$$

Therefore,  $l$  is an antiderivative of this last function. Using a trig substitution,  $x = r \sin \theta$ , it follows  $dx = r \cos \theta d\theta$  and so

$$\begin{aligned} \int \frac{r}{\sqrt{r^2 - x^2}} dx &= \int \frac{1}{\sqrt{1 - \sin^2 \theta}} r \cos \theta d\theta \\ &= r \int d\theta = r\theta + C \end{aligned}$$

Hence changing back to the variable  $x$  it follows  $l(x) = r \arcsin\left(\frac{x}{r}\right) + C$ . It only remains to find the constant. Plugging in  $x = 0$  this gives  $0 = C$  and  $l(x) = r \arcsin\left(\frac{x}{r}\right)$  so in particular, the length of the desired arc is

$$l\left(\frac{\sqrt{2}}{2}r\right) = r \arcsin\left(\frac{\sqrt{2}}{2}\right) = r \frac{\pi}{4}.$$

Note this gives the length of one eighth of the circle and so from this the length of the whole circle should be  $2r\pi$ . Here is another example

**Example 9.17.2** Find the length of the graph of  $y = x^2$  between  $x = 0$  and  $x = 1$ .

Here  $f'(x) = 2x$  and so the initial value problem to be solved is

$$\frac{dl}{dx} = \sqrt{1 + 4x^2}, \quad l(0) = 0.$$

Thus, getting the exact answer depends on finding

$$\int \sqrt{1 + 4x^2} dx.$$

Use the trig. substitution,  $2x = \tan u$  so  $dx = \frac{1}{2} (\sec^2 u) du$ . Therefore, making this substitution and using (9.10) on Page 194,

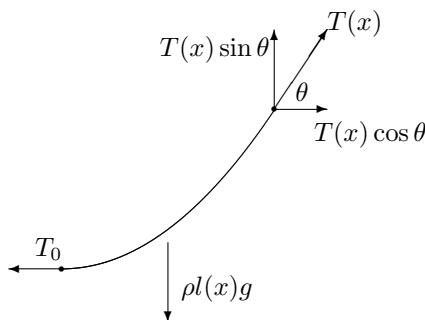
$$\begin{aligned} \int \sqrt{1 + 4x^2} dx &= \frac{1}{2} \int (\sec^3 u) du \\ &= \frac{1}{4} (\tan u) (\sec u) + \frac{1}{4} \ln |\sec u + \tan u| + C \\ &= \frac{1}{4} (2x) (\sqrt{1 + 4x^2}) + \frac{1}{4} \ln |2x + \sqrt{1 + 4x^2}| + C \end{aligned}$$

and since  $l(0) = 0$  it must be the case that  $C = 0$  and so the desired length is

$$l(1) = \frac{1}{2} \sqrt{5} + \frac{1}{4} \ln |2 + \sqrt{5}|$$

**Example 9.17.3** What is the equation of a hanging chain?

Consider the following picture of a portion of this chain.



In this picture,  $\rho$  denotes the density of the chain which is assumed to be constant and  $g$  is the acceleration due to gravity.  $T(x)$  and  $T_0$  represent the magnitude of the tension in the chain at  $x$  and at 0 respectively, as shown. Let the bottom of the chain be at the origin as shown. If this chain does not move, then all these forces acting on it must balance. In particular,

$$T(x) \sin \theta = l(x) \rho g, \quad T(x) \cos \theta = T_0.$$

Therefore, dividing these yields

$$\frac{\sin \theta}{\cos \theta} = l(x) \overbrace{\rho g / T_0}^{\equiv c}.$$

Now letting  $y(x)$  denote the  $y$  coordinate of the hanging chain corresponding to  $x$ ,

$$\frac{\sin \theta}{\cos \theta} = \tan \theta = y'(x).$$

Therefore, this yields

$$y'(x) = cl(x).$$

Now differentiating both sides of the differential equation,

$$y''(x) = cl'(x) = c\sqrt{1 + y'(x)^2}$$

and so

$$\frac{y''(x)}{\sqrt{1 + y'(x)^2}} = c.$$

Let  $z(x) = y'(x)$  so the above differential equation becomes

$$\frac{z'(x)}{\sqrt{1 + z^2}} = c.$$

Therefore,  $\int \frac{z'(x)}{\sqrt{1 + z^2}} dx = cx + d$ . Change the variable in the antiderivative letting  $u = z(x)$  and this yields

$$\begin{aligned} \int \frac{z'(x)}{\sqrt{1 + z^2}} dx &= \int \frac{du}{\sqrt{1 + u^2}} = \sinh^{-1}(u) + C \\ &= \sinh^{-1}(z(x)) + C. \end{aligned}$$

by (8.9) on Page 160. Therefore, combining the constants of integration,

$$\sinh^{-1}(y'(x)) = cx + d$$

and so

$$y'(x) = \sinh(cx + d).$$

Therefore,

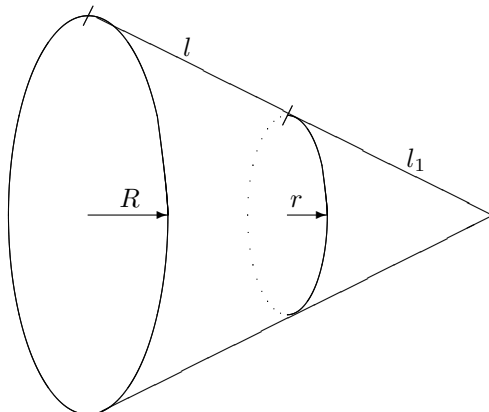
$$y(x) = \frac{1}{c} \cosh(cx + d) + k$$

where  $d$  and  $k$  are some constants and  $c = \rho g / T_0$ . Curves of this sort are called catenaries. Note these curves result from an assumption the only forces acting on the chain are as shown.

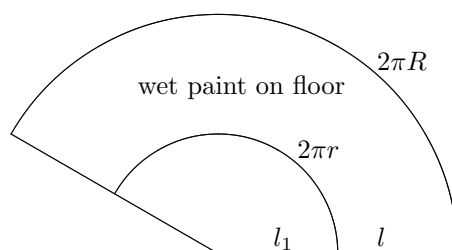
The problem of finding the surface area of a solid of revolution is closely related to that of finding the length of a graph. First consider the following picture of the frustum of a



cone in which it is desired to find the lateral surface area. In this picture, the frustum of the cone is the left part which has an  $l$  next to it and the lateral surface area is this part of the area of the cone.



To do this, imagine painting the sides and rolling the shape on the floor for exactly one revolution. The wet paint would make the following shape.



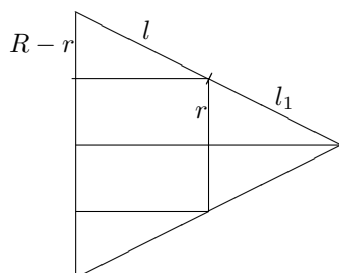
What would be the area of this wet paint? Its area would be the difference between the areas of the two sectors shown, one having radius  $l_1$  and the other having radius  $l + l_1$ . Both of these have the same central angle equal to

$$\frac{2\pi R}{2\pi(l + l_1)} 2\pi = \frac{2\pi R}{l + l_1}.$$

Therefore, by Theorem 3.8.2 on Page 66, this area is

$$(l + l_1)^2 \frac{\pi R}{(l + l_1)} - l_1^2 \frac{\pi R}{(l + l_1)} = \pi R l \frac{l + 2l_1}{l + l_1}$$

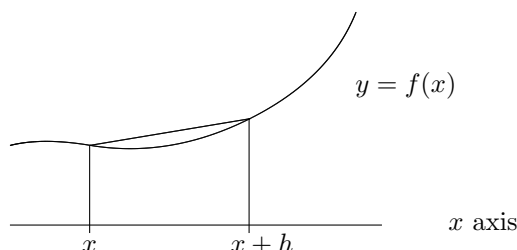
The view from the side is



and so by similar triangles,  $l_1 = lr / (R - r)$ . Therefore, substituting this into the above, the area of this frustum is

$$\pi R l \frac{l + 2 \left( \frac{lr}{R-r} \right)}{l + \left( \frac{lr}{R-r} \right)} = \pi l (R + r) = 2\pi l \left( \frac{R + r}{2} \right).$$

Now consider a function,  $f$ , defined on an interval,  $[a, b]$  and suppose it is desired to find the area of the surface which results when the graph of this function is revolved about the  $x$  axis. Consider the following picture of a piece of this graph.



Let  $A(x)$  denote the area which results from revolving the graph of the function restricted to  $[a, x]$ . Then from the above formula for the area of a frustum,

$$\frac{A(x+h) - A(x)}{h} \approx 2\pi \frac{1}{h} \sqrt{h^2 + (f(x+h) - f(x))^2} \left( \frac{f(x+h) + f(x)}{2} \right)$$

where  $\approx$  denotes that these are close to being equal and the approximation gets increasingly good as  $h \rightarrow 0$ . Therefore, rewriting this a little yields

$$\frac{A(x+h) - A(x)}{h} \approx 2\pi \sqrt{1 + \left( \frac{f(x+h) - f(x)}{h} \right)^2} \left( \frac{f(x+h) + f(x)}{2} \right)$$

Therefore, taking the limit as  $h \rightarrow 0$ , and using  $A(a) = 0$ , this yields the following initial value problem for  $A$  which can be used to find the area of a surface of revolution.

$$A'(x) = 2\pi f(x) \sqrt{1 + f'(x)^2}, \quad A(a) = 0.$$

**Example 9.17.4** Find the surface area of the surface obtained by revolving the function,  $y = r$  for  $x \in [a, b]$  about the  $x$  axis. Of course this is just the cylinder of radius  $r$  and height  $b - a$  so this area should equal  $2\pi r(b - a)$ . (Imagine painting it and rolling it on the floor and then taking the area of the rectangle which results.)

Using the above initial value problem, solve

$$A'(x) = 2\pi r \sqrt{1 + 0^2}, \quad A(a) = 0.$$

The solution is  $A(x) = 2\pi r(x - a)$ . Therefore,  $A(b) = 2\pi r(b - a)$  as expected.

**Example 9.17.5** Find the surface area of a sphere of radius  $r$ .

Here the function involved is  $f(x) = \sqrt{r^2 - x^2}$  for  $x \in [-r, r]$  and it is to be revolved about the  $x$  axis. In this case

$$f'(x) = \frac{-x}{\sqrt{r^2 - x^2}}$$

and so the initial value problem is of the form

$$A'(x) = 2\pi \sqrt{r^2 - x^2} \sqrt{1 + \frac{x^2}{r^2 - x^2}}, \quad A(-r) = 0$$

Thus, simplifying the above yields  $A'(x) = 2\pi r$  and so  $A(x) = 2\pi r x + C$  and since  $A(-r) = 0$ , it follows that  $C = 2\pi r^2$ . Therefore, the surface area equals  $A(r) = 2\pi r^2 + 2\pi r^2 = 4\pi r^2$ .

## 9.18 Exercises

1. In the hanging chain problem the picture and the derivation involved an assumption that at its lowest point, the chain was horizontal. Imagine lifting the end higher and higher and you will see this might not be the case in general. Can you modify the above derivation for the hanging chain to show that even in this case the chain will be in the form of a catenary?
2. Find the length of the graph of  $y = \ln(\cos x)$  for  $x \in [0, \pi/4]$ .
3. Find the length of the graph of  $y = x^{1/2} - \frac{x^{3/2}}{3}$  for  $x \in [0, 3]$ .
4. Find the length of  $y = \cosh(x)$  for  $x \in [0, 1]$ .
5. For  $a$  a positive real number, find the length of  $y = \frac{ax^2}{2} - \frac{1}{4a} \ln x$  for  $x \in [1, 2]$ . Of course your answer should depend on  $a$ .
6. The giant arch in St. Louis is in the form of an inverted catenary. Why?
7. The graph of the function,  $y = x^2$  for  $x \in [0, 1]$  is revolved about the  $x$  axis. Find the area of the surface of revolution.
8. The graph of the function,  $y = \sqrt{x}$  for  $x \in [0, 1]$  is revolved about the  $y$  axis. Find the area of the surface of revolution. **Hint:** Switch  $x$  and  $y$  and then use the previous problem.
9. The graph of the function,  $y = x^{1/2} - \frac{x^{3/2}}{3}$  is revolved about the  $x$  axis. Find the area of the surface of revolution if  $x \in [0, 2]$ .
10. The graph of the function,  $y = \cosh x$  for  $x \in [0, 1]$  is revolved about the  $x$  axis. Find the area of the surface of revolution.
11. The graph of the function,  $y = \sinh x$  for  $x \in [0, 1]$  is revolved about the  $x$  axis. Find the area of the surface of revolution.

12. The ellipse,  $\frac{x^2}{a^2} + \frac{y^2}{b^2} = 1$  is revolved about the  $x$  axis. Find the area of the surface of revolution.
13. Separable differential equations are those which can be written in the form

$$\frac{dy}{dx} = \frac{f(x)}{g(y)}.$$

The reason these are called separable is that if you formally cross multiply,

$$g(y) dy = f(x) dx$$

and the variables are “separated”. The  $x$  variables are on one side and the  $y$  variables are on the other. Such an equation was encountered in the hanging chain problem. Recall  $z'(x)/\sqrt{1+z^2} = c$  so  $\frac{dz}{dx} = c\sqrt{1+z^2}$ . Now show that if  $G'(y) = g(y)$  and  $F'(x) = f(x)$ , then if the equation,  $F(x) - G(y) = c$  specifies  $y$  as a differentiable function of  $x$ , then  $x \rightarrow y(x)$  solves the separable differential equation. For this reason the expression,  $F(x) - G(y) = c$  for  $c$  an arbitrary constant is regarded as an acceptable description of the solution of the differential equation. In short, this indicates that to find the solutions of the equation, it suffices to take  $G(y) \in \int g(y) dy$ ,  $F(x) \in \int f(x) dx$ , and obtain the solutions as  $F(x) - G(y) = c$  where  $c$  is a constant.

**Hint:** Use implicit differentiation or the chain rule.

14. Find the solution to the initial value problem,

$$y' = \frac{x}{y^2}, \quad y(0) = 1.$$

15. Find the solution to the initial value problem,

$$y' = 1 + y^2, \quad y(0) = 0.$$

## 9.19 The Equation $y' + ay = 0$

The homogeneous first order constant coefficient linear differential equation is a differential equation of the form

$$y' + ay = 0. \tag{9.15}$$

It is arguably the most important differential equation in existence. Generalizations of it include the entire subject of linear differential equations and even many of the most important partial differential equations occurring in applications.

Here is how to find the solutions to this equation. Multiply both sides of the equation by  $e^{at}$ . Then use the product and chain rules to verify that

$$e^{at}(y' + ay) = \frac{d}{dt}(e^{at}y) = 0.$$

Therefore, since the derivative of the function  $t \rightarrow e^{at}y(t)$  equals zero, it follows this function must equal some constant,  $C$ . Consequently,  $ye^{at} = C$  and so  $y(t) = Ce^{-at}$ . This shows that if there is a solution of the equation,  $y' + ay = 0$ , then it must be of the form  $Ce^{-at}$  for some constant,  $C$ . You should verify that every function of the form,  $y(t) = Ce^{-at}$  is a solution of the above differential equation, showing this yields all solutions. This proves the following theorem.

**Theorem 9.19.1** *The solutions to the equation,  $y' + ay = 0$  consist of all functions of the form,  $Ce^{-at}$  where  $C$  is some constant.*

**Exercise 9.19.2** *Radioactive substances decay in the following way. The rate of decay is proportional to the amount present. In other words, letting  $A(t)$  denote the amount of the radioactive substance at time  $t$ ,  $A(t)$  satisfies the following initial value problem.*

$$A'(t) = -k^2 A(t), \quad A(0) = A_0$$

where  $A_0$  is the initial amount of the substance. What is the solution to the initial value problem?

Write the differential equation as  $A'(t) + k^2 A(t) = 0$ . From Theorem 9.19.1 the solution is

$$A(t) = Ce^{-k^2 t}$$

and it only remains to find  $C$ . Letting  $t = 0$ , it follows  $A_0 = A(0) = C$ . Thus  $A(t) = A_0 \exp(-k^2 t)$ .

## 9.20 Exercises

1. For  $x$  sufficiently large, let  $f(x) = (\ln x)^{\ln x}$ . Find  $f'(x)$ . How big does  $x$  need to be in order for this to make sense. **Hint:** You should have  $\ln x > 0$ .
2. Verify that for any constant,  $C$ , the function,  $y(t) = Ce^{-at}$  solves the differential equation,  $y' + ay = 0$ .
3. Prove Theorem 9.19.1 by using Problem 13 on Page 220.
4. In Example 9.19.2 the half life is the time it takes for half of the original amount of the substance to decay. Find a formula for the half life assuming you know  $k^2$ .
5. There are ten grams of a radioactive substance which is allowed to decay for five years. At the end of the five years there are 9.5 grams of the substance left. Find the half life of the substance. **Hint:** Use the given information to find  $k^2$  and then use Problem 4.
6. Sometimes banks compound interest continuously. They do so by letting the amount in the account satisfy the initial value problem,

$$A'(t) = rA(t), \quad A(0) = A_0$$

where here  $A(t)$  is the amount at time  $t$  measured in years,  $r$  is the interest rate per year, and  $A_0$  is the initial amount. Find  $A(t)$  explicitly. If \$100 is placed in an account which is compounded continuously at 6% per year, how many years will it be before there is \$200 in the account? **Hint:** In this case,  $r = .06$ .

7. A population is growing at the rate of 11% per year. This means it satisfies the differential equation,  $A' = .11A$ . Find the time it takes for the population to double.
8. A substance is decaying at the rate of .01 per year. Find the half life of the substance.
9. The half live of a substance is 40 years. Find the rate of decay of the substance.
10. A sample of 4 grams of a radioactive substance has decayed to 2 grams in 5 days. Find the half life of the substance. Give your answer in terms of logarithms.

11. \$1000 is deposited in an account that earns interest at the rate of 5% per year compounded continuously. How much will be in the account after 10 years?
12. A sample of one ounce of water from a water supply is cultured in a Petri dish. After four hours, there are 3000.0 bacteria present and after six hours there are 7000 bacteria present. How many were present in the original sample?
13. Carbon 14 is a radioactive isotope of carbon and it is produced at a more or less constant rate in the earth's atmosphere by radiation from the sun. It also decays at a rate proportional to the amount present. Show this implies that, assuming this has been going on for billions of years, it is reasonable to assume the amount of Carbon 14 in the atmosphere is essentially constant over time. **Hint:** By assumption, if  $A$  is the amount of Carbon 14 in the atmosphere,  $\frac{dA}{dt} = -k^2A + r$  where  $k^2$  is the constant of decay described above and  $r$  is the constant rate of production. Now differentiate both sides and show  $A'(t) = Ce^{-k^2t}$ . Conclude  $A(t) = D - (C/k^2)e^{-k^2t}$ . What happens to this over a long period of time?
14. The method of carbon dating is based on the result of Problem 13. When an animal or plant is alive it absorbs carbon from the atmosphere and so when it is living it has a known percentage of carbon 14. When it dies, it quits absorbing carbon and the carbon 14 begins to decay. After some time,  $t$ , the amount of carbon 14 can be measured and on this basis, the time since the death of the animal or plant can be estimated. Given the half life of carbon 14 is 5730 years and the amount of carbon 14 in a mummy is .3 what it was at the time of death, how long has it been since the mummy was alive? (To see how to come up with this figure for the half life, see Problem 5. You do experiments and take measurements over a smaller period of time.)
15. The half life of carbon 14 is known to be 5730 years. A certain tree stump is known to be 4000 years old. What percentage of the original carbon 14 should it contain?
16. One model for population growth is to assume the rate of growth is proportional to both the population and the difference between some maximum sustainable population and the population. The reason for this is that when the population gets large enough, there begin to be insufficient resources. Thus  $\frac{dA}{dt} = kA(M - A)$ , where  $k$  and  $M$  are positive constants. Show this is a separable differential equation and its solutions are of the form

$$A(t) = \frac{M}{1 + Ce^{-kMt}}$$

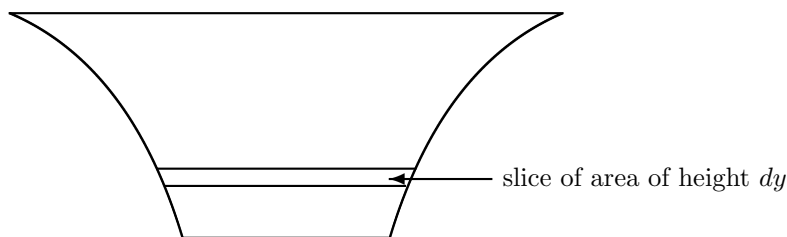
where  $C$  is a constant. Given three measurements of population at three equally spaced times, show how to predict the maximum sustainable population<sup>1</sup>.

## 9.21 Force On A Dam And Work Of A Pump

Imagine you are a fish swimming in a lake behind a dam and you are interested in the total force acting on the dam. The following picture is what you would see.

---

<sup>1</sup>This has been done with the earth's population and the maximum sustainable population has been exceeded. Therefore, the model is far too simplistic for human population growth. However, it would work somewhat better for predicting the growth of things like bacteria.



The reason you would be interested in that long thin slice of area having essentially the same depth, say at  $y$  feet is because the pressure in the water at that depth is constant and equals  $62.5y$  pounds per square foot<sup>2</sup>. Therefore, the total force the water exerts on the long thin slice is

$$dF = 62.5yL(y) dy$$

where  $L(y)$  denotes the length of the slice. Therefore, the total force on the dam up to depth  $y$  is obtained as a solution to the initial value problem,

$$\frac{dF}{dy} = 62.5yL(y), \quad F(0) = 0.$$

**Example 9.21.1** Suppose the width of a dam at depth  $y$  feet equals  $L(y) = 1000 - y$  and its depth is 500 feet. Find the total force in pounds exerted on the dam.

From the above, this is obtained as the solution to the initial value problem

$$\frac{dF}{dy} = 62.5y(1000 - y), \quad F(0) = 0$$

which is  $F(y) = -20.83y^3 + 31250y^2$ . The total force on the dam would be

$$F(500) = -20.83(500)^3 + 31250(500)^2 = 5,208,750,000.0$$

pounds. In tons this is 2,604,375. That is a lot of force.

Now suppose you are pumping water from a tank of depth  $d$  to a height of  $H$  feet above the top of the water in the tank. Suppose also that at depth  $y$  below the surface, the area of a cross section having constant depth is  $A(y)$ . The total weight of a slice of water having thickness  $dy$  at this depth is  $62.5A(y)dy$  and the pump needs to lift this weight a distance of  $y + H$  feet. Therefore, the work done is  $dW = (y + H)62.5A(y)dy$ . An initial value problem for the work done to pump the water down to a depth of  $y$  feet would be

$$\frac{dW}{dy} = (y + H)62.5A(y), \quad W(0) = 0.$$

The reason for the initial condition is that the pump has done no work to pump no water. If the weight of the fluid per cubic foot were different than 62.5 you would do the same things but replace the number.

<sup>2</sup>Later on a nice result on hydrostatic pressure will be presented which will verify this assertion. Here 62.5 is the weight in pounds of a cubic foot of water. If you like, think of a column of water of height  $y$  having base area equal to 1 square foot. Then the total force acting on this base area would be  $62.5 \times y$  pounds.

**Example 9.21.2** *A spherical storage tank sitting on the ground having radius 30 feet is half filled with a fluid which weighs 50 pounds per cubic foot. How much work is done to pump this fluid to a height of 100 feet?*

Letting  $r$  denote the radius of a cross section  $y$  feet below the level of the fluid,  $r^2 + y^2 = 900$ . Therefore,

$$r = \sqrt{900 - y^2}.$$

It follows the area of the cross section at depth  $y$  is  $\pi(900 - y^2)$ . Here  $H = 70$  and so the initial value problem to solve is

$$\frac{dW}{dy} = (y + 70) 50\pi (900 - y^2), \quad W(0) = 0.$$

Therefore,  $W(y) = 50\pi \left(-\frac{1}{4}y^4 - \frac{70}{3}y^3 + 450y^2 + 63\,000y\right)$  and the total work in foot pounds equals

$$W(30) = 50\pi \left(-\frac{1}{4}(30)^4 - \frac{70}{3}(30)^3 + 450(30)^2 + 63\,000(30)\right) = 73,125,000\pi$$

## 9.22 Exercises

1. A cylindrical storage tank having radius 20 feet and length 40 feet is filled with a fluid which weighs 50 pounds per cubic foot. This tank is lying on its side on the ground. Find the total force acting on the ends of the tank by the fluid.
2. Suppose the tank in the above problem is filled to a depth of 8 feet. Find the work needed to pump the fluid to a height of 50 feet.
3. A conical hole is filled with water. If the depth of the hole is 20 feet and the radius of the hole is 10 feet, how much work is needed to pump the water to a height of 10 feet above the ground?
4. A dam 500 feet high has a width at depth  $y$  equal to  $4000 - 2y$  feet. What is the total force on the dam if it is filled?
5. When the bucket is filled with water it weighs 30 pounds and when empty it weighs 2 pounds and the person on top of a 100 foot building exerts a constant force of 40 pounds. The bucket is full at the bottom but leaks at the rate of .1 cubic feet per second. How much work does the person on the top of the building do in lifting the bucket to the top? Will the bucket be empty when it reaches the top? You can use Newton's law that force equals mass times acceleration.
6. In the situation of the above problem, suppose the person on the top maintains a constant velocity of 1 foot per second. How much work does he do and is the bucket empty when it reaches the top?
7. A silo is 10 feet in diameter and at a height of 30 feet there is a hemispherical top. The silage weighs 10 pounds per cubic foot. How much work was done in filling it to the very top?
8. A cylindrical storage tank having radius 10 feet is filled with water to a depth of 20 feet. If the storage tank stands upright on its circular base, what is the total force the water exerts on the sides of the tank? **Hint:** The pressure in the water at depth  $y$  is  $62.5y$  pounds per square foot.



9. A spherical storage tank having radius 10 feet is filled with water. What is the total force the water exerts on the storage tank? **Hint:** The pressure in the water at depth  $y$  is  $62.5y$  consider the area corresponding to a slice at height  $y$ . This is a surface of revolution and you know how to deal with these.



# The Integral

The integral originated in attempts to find areas of various shapes and the ideas involved in finding integrals are much older than the ideas related to finding derivatives. In fact, Archimedes<sup>1</sup> was finding areas of various curved shapes about 250 B.C. The integral is needed to remove some of the mathematical loose ends and also to enable the study of more general problems involving differential equations. It will also be useful for formulating other physical models. The technique used for finding the area of a circular segment presented early in the book was essentially that employed by Archimedes and contains the essential ideas for the integral. The main difference is that here the triangles will be replaced with rectangles. You may be wondering what the fuss is about. Areas have already been found as solutions of differential equations. However, there is a profound difference between what is about to be presented and what has just been done. It is related to the fundamental mathematical question of existence. As an illustration, consider the problem of finding the area between  $y = e^{x^2}$  and the  $x$  axis for  $x \in [0, 1]$ . As pointed out earlier, the area is obtained as a solution to the initial value problem,

$$A'(x) = e^{x^2}, \quad A(0) = 0.$$

So what is the solution to this initial value problem? By Theorem 9.1.1 there is at most one solution, but what is the solution? Does it even exist? More generally, for which functions,  $f$  does there exist a solution to the initial value problem,  $y'(x) = f(x), y(0) = y_0$ ? These questions are typical of mathematics. There are usually two aspects to a mathematical concept. One is the question of existence and the other is how to find that which exists. The two questions are often very different and one can have a good understanding of one without having any idea how to go about considering the other. However, both are absolutely essential. In the preceding chapter the only thing considered was the second question.

## 10.1 Upper And Lower Sums

The Riemann integral pertains to bounded functions which are defined on a bounded interval. Let  $[a, b]$  be a closed interval. A set of points in  $[a, b]$ ,  $\{x_0, \dots, x_n\}$  is a partition if

$$a = x_0 < x_1 < \dots < x_n = b.$$

---

<sup>1</sup>Archimedes 287-212 B.C. found areas of curved regions by stuffing them with simple shapes which he knew the area of and taking a limit. He also made fundamental contributions to physics. The story is told about how he determined that a gold smith had cheated the king by giving him a crown which was not solid gold as had been claimed. He did this by finding the amount of water displaced by the crown and comparing with the amount of water it should have displaced if it had been solid gold.

Such partitions are denoted by  $P$  or  $Q$ . For  $f$  a bounded function defined on  $[a, b]$ , let

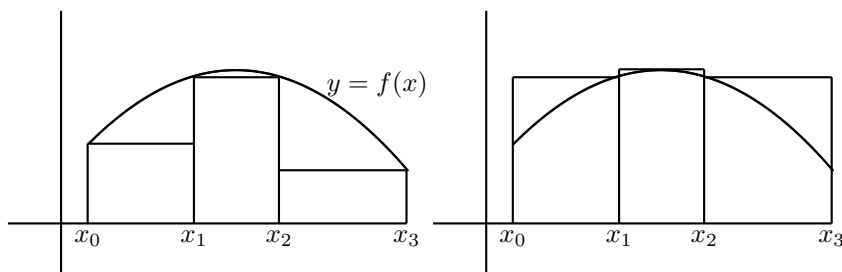
$$M_i(f) \equiv \sup\{f(x) : x \in [x_{i-1}, x_i]\},$$

$$m_i(f) \equiv \inf\{f(x) : x \in [x_{i-1}, x_i]\}.$$

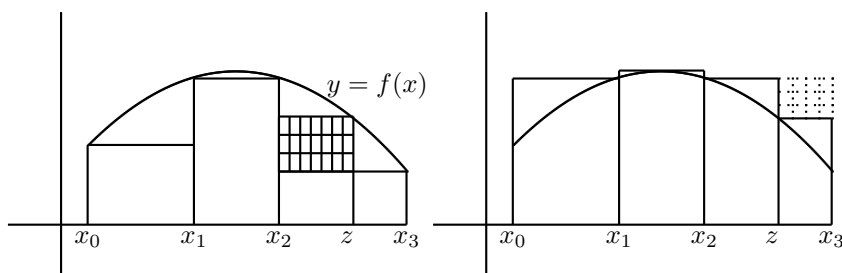
Also let  $\Delta x_i \equiv x_i - x_{i-1}$ . Then define upper and lower sums as

$$U(f, P) \equiv \sum_{i=1}^n M_i(f) \Delta x_i \text{ and } L(f, P) \equiv \sum_{i=1}^n m_i(f) \Delta x_i$$

respectively. The numbers,  $M_i(f)$  and  $m_i(f)$ , are well defined real numbers because  $f$  is assumed to be bounded and  $\mathbb{R}$  is complete. Thus the set  $S = \{f(x) : x \in [x_{i-1}, x_i]\}$  is bounded above and below. In the following picture, the sum of the areas of the rectangles in the picture on the left is a lower sum for the function in the picture and the sum of the areas of the rectangles in the picture on the right is an upper sum for the same function which uses the same partition.



What happens when you add in more points in a partition? The following pictures illustrate in the context of the above example. In this example a single additional point, labeled  $z$  has been added in.



Note how the lower sum got larger by the amount of the area in the shaded rectangle and the upper sum got smaller by the amount in the rectangle shaded by dots. In general this is the way it works and this is shown in the following lemma.

**Lemma 10.1.1** *If  $P \subseteq Q$  then*

$$U(f, Q) \leq U(f, P), \text{ and } L(f, P) \leq L(f, Q).$$

**Proof:** This is verified by adding in one point at a time. Thus let  $P = \{x_0, \dots, x_n\}$  and let  $Q = \{x_0, \dots, x_k, y, x_{k+1}, \dots, x_n\}$ . Thus exactly one point,  $y$ , is added between  $x_k$

and  $x_{k+1}$ . Now the term in the upper sum which corresponds to the interval  $[x_k, x_{k+1}]$  in  $U(f, P)$  is

$$\sup \{f(x) : x \in [x_k, x_{k+1}]\} (x_{k+1} - x_k) \quad (10.1)$$

and the term which corresponds to the interval  $[x_k, x_{k+1}]$  in  $U(f, Q)$  is

$$\sup \{f(x) : x \in [x_k, y]\} (y - x_k) + \sup \{f(x) : x \in [y, x_{k+1}]\} (x_{k+1} - y) \quad (10.2)$$

$$\equiv M_1 (y - x_k) + M_2 (x_{k+1} - y) \quad (10.3)$$

All the other terms in the two sums coincide. Now  $\sup \{f(x) : x \in [x_k, x_{k+1}]\} \geq \max(M_1, M_2)$  and so the expression in (10.2) is no larger than

$$\begin{aligned} & \sup \{f(x) : x \in [x_k, x_{k+1}]\} (x_{k+1} - y) + \sup \{f(x) : x \in [x_k, x_{k+1}]\} (y - x_k) \\ &= \sup \{f(x) : x \in [x_k, x_{k+1}]\} (x_{k+1} - x_k), \end{aligned}$$

the term corresponding to the interval,  $[x_k, x_{k+1}]$  and  $U(f, P)$ . This proves the first part of the lemma pertaining to upper sums because if  $Q \supseteq P$ , one can obtain  $Q$  from  $P$  by adding in one point at a time and each time a point is added, the corresponding upper sum either gets smaller or stays the same. The second part is similar and is left as an exercise.

**Lemma 10.1.2** *If  $P$  and  $Q$  are two partitions, then*

$$L(f, P) \leq U(f, Q).$$

**Proof:** By Lemma 10.1.1,

$$L(f, P) \leq L(f, P \cup Q) \leq U(f, P \cup Q) \leq U(f, Q).$$

**Definition 10.1.3**

$$\bar{I} \equiv \inf \{U(f, Q) \text{ where } Q \text{ is a partition}\}$$

$$\underline{I} \equiv \sup \{L(f, P) \text{ where } P \text{ is a partition}\}.$$

Note that  $\underline{I}$  and  $\bar{I}$  are well defined real numbers.

**Theorem 10.1.4**  $\underline{I} \leq \bar{I}$ .

**Proof:** From Lemma 10.1.2,

$$\underline{I} = \sup \{L(f, P) \text{ where } P \text{ is a partition}\} \leq U(f, Q)$$

because  $U(f, Q)$  is an upper bound to the set of all lower sums and so it is no smaller than the least upper bound. Therefore, since  $Q$  is arbitrary,

$$\begin{aligned} \underline{I} &= \sup \{L(f, P) \text{ where } P \text{ is a partition}\} \\ &\leq \inf \{U(f, Q) \text{ where } Q \text{ is a partition}\} \equiv \bar{I} \end{aligned}$$

where the inequality holds because it was just shown that  $\underline{I}$  is a lower bound to the set of all upper sums and so it is no larger than the greatest lower bound of this set. This proves the theorem.

**Definition 10.1.5** *A bounded function  $f$  is Riemann integrable, written as*

$$f \in R([a, b])$$

*if*

$$\underline{I} = \bar{I}$$

*and in this case,*

$$\int_a^b f(x) dx \equiv \underline{I} = \bar{I}.$$

Thus, in words, the Riemann integral is the unique number which lies between all upper sums and all lower sums if there is such a unique number.

Recall Proposition 2.14.3. It is stated here for ease of reference.

**Proposition 10.1.6** *Let  $S$  be a nonempty set and suppose  $\sup(S)$  exists. Then for every  $\delta > 0$ ,*

$$S \cap (\sup(S) - \delta, \sup(S)] \neq \emptyset.$$

*If  $\inf(S)$  exists, then for every  $\delta > 0$ ,*

$$S \cap [\inf(S), \inf(S) + \delta) \neq \emptyset.$$

This proposition implies the following theorem which is used to determine the question of Riemann integrability.

**Theorem 10.1.7** *A bounded function  $f$  is Riemann integrable if and only if for all  $\varepsilon > 0$ , there exists a partition  $P$  such that*

$$U(f, P) - L(f, P) < \varepsilon. \quad (10.4)$$

**Proof:** First assume  $f$  is Riemann integrable. Then let  $P$  and  $Q$  be two partitions such that

$$U(f, Q) < \bar{I} + \varepsilon/2, \quad L(f, P) > \underline{I} - \varepsilon/2.$$

Then since  $\underline{I} = \bar{I}$ ,

$$U(f, Q \cup P) - L(f, P \cup Q) \leq U(f, Q) - L(f, P) < \bar{I} + \varepsilon/2 - (\underline{I} - \varepsilon/2) = \varepsilon.$$

Now suppose that for all  $\varepsilon > 0$  there exists a partition such that (10.4) holds. Then for given  $\varepsilon$  and partition  $P$  corresponding to  $\varepsilon$

$$\bar{I} - \underline{I} \leq U(f, P) - L(f, P) \leq \varepsilon.$$

Since  $\varepsilon$  is arbitrary, this shows  $\underline{I} = \bar{I}$  and this proves the theorem.

The condition described in the theorem is called the Riemann criterion.

Not all bounded functions are Riemann integrable. For example, let

$$f(x) \equiv \begin{cases} 1 & \text{if } x \in \mathbb{Q} \\ 0 & \text{if } x \in \mathbb{R} \setminus \mathbb{Q} \end{cases} \quad (10.5)$$

Then if  $[a, b] = [0, 1]$  all upper sums for  $f$  equal 1 while all lower sums for  $f$  equal 0. Therefore the Riemann criterion is violated for  $\varepsilon = 1/2$ .

## 10.2 Exercises

1. Prove the second half of Lemma 10.1.1 about lower sums.
2. Verify that for  $f$  given in (10.5), the lower sums on the interval  $[0, 1]$  are all equal to zero while the upper sums are all equal to one.
3. Let  $f(x) = 1 + x^2$  for  $x \in [-1, 3]$  and let  $P = \{-1, -\frac{1}{3}, 0, \frac{1}{2}, 1, 2\}$ . Find  $U(f, P)$  and  $L(f, P)$ .

4. Show that if  $f \in R([a, b])$ , there exists a partition,  $\{x_0, \dots, x_n\}$  such that for any  $z_k \in [x_k, x_{k+1}]$ ,

$$\left| \int_a^b f(x) dx - \sum_{k=1}^n f(z_k)(x_k - x_{k-1}) \right| < \varepsilon$$

This sum,  $\sum_{k=1}^n f(z_k)(x_k - x_{k-1})$ , is called a Riemann sum and this exercise shows that the integral can always be approximated by a Riemann sum.

5. Let  $P = \{1, 1\frac{1}{4}, 1\frac{1}{2}, 1\frac{3}{4}, 2\}$ . Find upper and lower sums for the function,  $f(x) = \frac{1}{x}$  using this partition. What does this tell you about  $\ln(2)$ ?
6. If  $f \in R([a, b])$  and  $f$  is changed at finitely many points, show the new function is also in  $R([a, b])$ .
7. Define a “left sum” as

$$\sum_{k=1}^n f(x_{k-1})(x_k - x_{k-1})$$

and a “right sum”,

$$\sum_{k=1}^n f(x_k)(x_k - x_{k-1}).$$

Also suppose that all partitions have the property that  $x_k - x_{k-1}$  equals a constant,  $(b - a)/n$  so the points in the partition are equally spaced, and define the integral to be the number these right and left sums get close to as  $n$  gets larger and larger. Show that for  $f$  given in (10.5),  $\int_0^x f(t) dt = 1$  if  $x$  is rational and  $\int_0^x f(t) dt = 0$  if  $x$  is irrational. It turns out that the correct answer should always equal zero for that function, regardless of whether  $x$  is rational. This is shown in more advanced courses when the Lebesgue integral is studied. This illustrates why this method of defining the integral in terms of left and right sums is total nonsense.

## 10.3 Functions Of Riemann Integrable Functions

It is often necessary to consider functions of Riemann integrable functions and a natural question is whether these are Riemann integrable. The following theorem gives a partial answer to this question. This is not the most general theorem which will relate to this question but it will be enough for the needs of this book.

**Theorem 10.3.1** *Let  $f, g$  be bounded functions and let  $f([a, b]) \subseteq [c_1, d_1]$  and  $g([a, b]) \subseteq [c_2, d_2]$ . Let  $H : [c_1, d_1] \times [c_2, d_2] \rightarrow \mathbb{R}$  satisfy,*

$$|H(a_1, b_1) - H(a_2, b_2)| \leq K[|a_1 - a_2| + |b_1 - b_2|]$$

*for some constant  $K$ . Then if  $f, g \in R([a, b])$  it follows that  $H \circ (f, g) \in R([a, b])$ .*

**Proof:** In the following claim,  $M_i(h)$  and  $m_i(h)$  have the meanings assigned above with respect to some partition of  $[a, b]$  for the function,  $h$ .

**Claim:** The following inequality holds.

$$\begin{aligned} & |M_i(H \circ (f, g)) - m_i(H \circ (f, g))| \leq \\ & K[|M_i(f) - m_i(f)| + |M_i(g) - m_i(g)|]. \end{aligned}$$

**Proof of the claim:** By the above proposition, there exist  $x_1, x_2 \in [x_{i-1}, x_i]$  be such that

$$H(f(x_1), g(x_1)) + \eta > M_i(H \circ (f, g)),$$

and

$$H(f(x_2), g(x_2)) - \eta < m_i(H \circ (f, g)).$$

Then

$$\begin{aligned} & |M_i(H \circ (f, g)) - m_i(H \circ (f, g))| \\ & < 2\eta + |H(f(x_1), g(x_1)) - H(f(x_2), g(x_2))| \\ & < 2\eta + K[|f(x_1) - f(x_2)| + |g(x_1) - g(x_2)|] \\ & \leq 2\eta + K[|M_i(f) - m_i(f)| + |M_i(g) - m_i(g)|]. \end{aligned}$$

Since  $\eta > 0$  is arbitrary, this proves the claim.

Now continuing with the proof of the theorem, let  $P$  be such that

$$\sum_{i=1}^n (M_i(f) - m_i(f)) \Delta x_i < \frac{\varepsilon}{2K}, \quad \sum_{i=1}^n (M_i(g) - m_i(g)) \Delta x_i < \frac{\varepsilon}{2K}.$$

Then from the claim,

$$\begin{aligned} & \sum_{i=1}^n (M_i(H \circ (f, g)) - m_i(H \circ (f, g))) \Delta x_i \\ & < \sum_{i=1}^n K[|M_i(f) - m_i(f)| + |M_i(g) - m_i(g)|] \Delta x_i < \varepsilon. \end{aligned}$$

Since  $\varepsilon > 0$  is arbitrary, this shows  $H \circ (f, g)$  satisfies the Riemann criterion and hence  $H \circ (f, g)$  is Riemann integrable as claimed. This proves the theorem.

This theorem implies that if  $f, g$  are Riemann integrable, then so is  $af + bg, |f|, f^2$ , along with infinitely many other such continuous combinations of Riemann integrable functions. For example, to see that  $|f|$  is Riemann integrable, let  $H(a, b) = |a|$ . Clearly this function satisfies the conditions of the above theorem and so  $|f| = H(f, f) \in R([a, b])$  as claimed. The following theorem gives an example of many functions which are Riemann integrable.

**Theorem 10.3.2** Let  $f : [a, b] \rightarrow \mathbb{R}$  be either increasing or decreasing on  $[a, b]$ . Then  $f \in R([a, b])$ .

**Proof:** Let  $\varepsilon > 0$  be given and let

$$x_i = a + i \left( \frac{b-a}{n} \right), \quad i = 0, \dots, n.$$

Then since  $f$  is increasing,

$$\begin{aligned} U(f, P) - L(f, P) &= \sum_{i=1}^n (f(x_i) - f(x_{i-1})) \left( \frac{b-a}{n} \right) \\ &= (f(b) - f(a)) \left( \frac{b-a}{n} \right) < \varepsilon \end{aligned}$$

whenever  $n$  is large enough. Thus the Riemann criterion is satisfied and so the function is Riemann integrable. The proof for decreasing  $f$  is similar.



**Corollary 10.3.3** *Let  $[a, b]$  be a bounded closed interval and let  $\phi : [a, b] \rightarrow \mathbb{R}$  be Lipschitz continuous. Then  $\phi \in R([a, b])$ . Recall that a function,  $\phi$ , is Lipschitz continuous if there is a constant,  $K$ , such that for all  $x, y$ ,*

$$|\phi(x) - \phi(y)| < K|x - y|.$$

**Proof:** Let  $f(x) = x$ . Then by Theorem 10.3.2,  $f$  is Riemann integrable. Let  $H(a, b) \equiv \phi(a)$ . Then by Theorem 10.3.1  $H \circ (f, f) = \phi \circ f = \phi$  is also Riemann integrable. This proves the corollary.

In fact, it is enough to assume  $\phi$  is continuous, although this is harder. This is the content of the next theorem which is where the difficult theorems about continuity and uniform continuity are used.

**Theorem 10.3.4** *Suppose  $f : [a, b] \rightarrow \mathbb{R}$  is continuous. Then  $f \in R([a, b])$ .*

**Proof:** By Corollary 5.13.5 on Page 114,  $f$  is uniformly continuous on  $[a, b]$ . Therefore, if  $\varepsilon > 0$  is given, there exists a  $\delta > 0$  such that if  $|x_i - x_{i-1}| < \delta$ , then  $M_i - m_i < \frac{\varepsilon}{b-a}$ . Let

$$P \equiv \{x_0, \dots, x_n\}$$

be a partition with  $|x_i - x_{i-1}| < \delta$ . Then

$$U(f, P) - L(f, P) < \sum_{i=1}^n (M_i - m_i)(x_i - x_{i-1}) < \frac{\varepsilon}{b-a}(b-a) = \varepsilon.$$

By the Riemann criterion,  $f \in R([a, b])$ . This proves the theorem.

## 10.4 Properties Of The Integral

The integral has many important algebraic properties. First here is a simple lemma.

**Lemma 10.4.1** *Let  $S$  be a nonempty set which is bounded above and below. Then if  $-S \equiv \{-x : x \in S\}$ ,*

$$\sup(-S) = -\inf(S) \tag{10.6}$$

*and*

$$\inf(-S) = -\sup(S). \tag{10.7}$$

**Proof:** Consider (10.6). Let  $x \in S$ . Then  $-x \leq \sup(-S)$  and so  $x \geq -\sup(-S)$ . It follows that  $-\sup(-S)$  is a lower bound for  $S$  and therefore,  $-\sup(-S) \leq \inf(S)$ . This implies  $\sup(-S) \geq -\inf(S)$ . Now let  $-x \in -S$ . Then  $x \in S$  and so  $x \geq \inf(S)$  which implies  $-x \leq -\inf(S)$ . Therefore,  $-\inf(S)$  is an upper bound for  $-S$  and so  $-\inf(S) \geq \sup(-S)$ . This shows (10.6). Formula (10.7) is similar and is left as an exercise.

In particular, the above lemma implies that for  $M_i(f)$  and  $m_i(f)$  defined above  $M_i(-f) = -m_i(f)$ , and  $m_i(-f) = -M_i(f)$ .

**Lemma 10.4.2** *If  $f \in R([a, b])$  then  $-f \in R([a, b])$  and*

$$-\int_a^b f(x) dx = \int_a^b -f(x) dx.$$

**Proof:** The first part of the conclusion of this lemma follows from Theorem 10.3.2 since the function  $\phi(y) \equiv -y$  is Lipschitz continuous. Now choose  $P$  such that

$$\int_a^b -f(x) dx - L(-f, P) < \varepsilon.$$

Then since  $m_i(-f) = -M_i(f)$ ,

$$\varepsilon > \int_a^b -f(x) dx - \sum_{i=1}^n m_i(-f) \Delta x_i = \int_a^b -f(x) dx + \sum_{i=1}^n M_i(f) \Delta x_i$$

which implies

$$\varepsilon > \int_a^b -f(x) dx + \sum_{i=1}^n M_i(f) \Delta x_i \geq \int_a^b -f(x) dx + \int_a^b f(x) dx.$$

Thus, since  $\varepsilon$  is arbitrary,

$$\int_a^b -f(x) dx \leq - \int_a^b f(x) dx$$

whenever  $f \in R([a, b])$ . It follows

$$\int_a^b -f(x) dx \leq - \int_a^b f(x) dx = - \int_a^b -(-f(x)) dx \leq \int_a^b -f(x) dx$$

and this proves the lemma.

**Theorem 10.4.3** *The integral is linear,*

$$\int_a^b (\alpha f + \beta g)(x) dx = \alpha \int_a^b f(x) dx + \beta \int_a^b g(x) dx.$$

whenever  $f, g \in R([a, b])$  and  $\alpha, \beta \in \mathbb{R}$ .

**Proof:** First note that by Theorem 10.3.1,  $\alpha f + \beta g \in R([a, b])$ . To begin with, consider the claim that if  $f, g \in R([a, b])$  then

$$\int_a^b (f + g)(x) dx = \int_a^b f(x) dx + \int_a^b g(x) dx. \quad (10.8)$$

Let  $P_1, Q_1$  be such that

$$U(f, Q_1) - L(f, Q_1) < \varepsilon/2, \quad U(g, P_1) - L(g, P_1) < \varepsilon/2.$$

Then letting  $P \equiv P_1 \cup Q_1$ , Lemma 10.1.1 implies

$$U(f, P) - L(f, P) < \varepsilon/2, \quad \text{and} \quad U(g, P) - L(g, P) < \varepsilon/2.$$

Next note that

$$m_i(f + g) \geq m_i(f) + m_i(g), \quad M_i(f + g) \leq M_i(f) + M_i(g).$$

Therefore,

$$L(g + f, P) \geq L(f, P) + L(g, P), \quad U(g + f, P) \leq U(f, P) + U(g, P).$$

For this partition,

$$\begin{aligned}\int_a^b (f+g)(x) dx &\in [L(f+g, P), U(f+g, P)] \\ &\subseteq [L(f, P) + L(g, P), U(f, P) + U(g, P)]\end{aligned}$$

and

$$\int_a^b f(x) dx + \int_a^b g(x) dx \in [L(f, P) + L(g, P), U(f, P) + U(g, P)].$$

Therefore,

$$\begin{aligned}\left| \int_a^b (f+g)(x) dx - \left( \int_a^b f(x) dx + \int_a^b g(x) dx \right) \right| &\leq \\ U(f, P) + U(g, P) - (L(f, P) + L(g, P)) &< \varepsilon/2 + \varepsilon/2 = \varepsilon.\end{aligned}$$

This proves (10.8) since  $\varepsilon$  is arbitrary.

It remains to show that

$$\alpha \int_a^b f(x) dx = \int_a^b \alpha f(x) dx.$$

Suppose first that  $\alpha \geq 0$ . Then

$$\begin{aligned}\int_a^b \alpha f(x) dx &\equiv \sup\{L(\alpha f, P) : P \text{ is a partition}\} = \\ \alpha \sup\{L(f, P) : P \text{ is a partition}\} &\equiv \alpha \int_a^b f(x) dx.\end{aligned}$$

If  $\alpha < 0$ , then this and Lemma 10.4.2 imply

$$\begin{aligned}\int_a^b \alpha f(x) dx &= \int_a^b (-\alpha)(-f(x)) dx \\ &= (-\alpha) \int_a^b (-f(x)) dx = \alpha \int_a^b f(x) dx.\end{aligned}$$

This proves the theorem.

**Theorem 10.4.4** *If  $f \in R([a, b])$  and  $f \in R([b, c])$ , then  $f \in R([a, c])$  and*

$$\int_a^c f(x) dx = \int_a^b f(x) dx + \int_b^c f(x) dx. \quad (10.9)$$

**Proof:** Let  $P_1$  be a partition of  $[a, b]$  and  $P_2$  be a partition of  $[b, c]$  such that

$$U(f, P_i) - L(f, P_i) < \varepsilon/2, \quad i = 1, 2.$$

Let  $P \equiv P_1 \cup P_2$ . Then  $P$  is a partition of  $[a, c]$  and

$$\begin{aligned}U(f, P) - L(f, P) &= U(f, P_1) - L(f, P_1) + U(f, P_2) - L(f, P_2) < \varepsilon/2 + \varepsilon/2 = \varepsilon.\end{aligned} \quad (10.10)$$

Thus,  $f \in R([a, c])$  by the Riemann criterion and also for this partition,

$$\begin{aligned} \int_a^b f(x) dx + \int_b^c f(x) dx &\in [L(f, P_1) + L(f, P_2), U(f, P_1) + U(f, P_2)] \\ &= [L(f, P), U(f, P)] \end{aligned}$$

and

$$\int_a^c f(x) dx \in [L(f, P), U(f, P)].$$

Hence by (10.10),

$$\left| \int_a^c f(x) dx - \left( \int_a^b f(x) dx + \int_b^c f(x) dx \right) \right| < U(f, P) - L(f, P) < \varepsilon$$

which shows that since  $\varepsilon$  is arbitrary, (10.9) holds. This proves the theorem.

**Corollary 10.4.5** *Let  $[a, b]$  be a closed and bounded interval and suppose that*

$$a = y_1 < y_2 < \cdots < y_l = b$$

*and that  $f$  is a bounded function defined on  $[a, b]$  which has the property that  $f$  is either increasing on  $[y_j, y_{j+1}]$  or decreasing on  $[y_j, y_{j+1}]$  for  $j = 1, \dots, l-1$ . Then  $f \in R([a, b])$ .*

**Proof:** This follows from Theorem 10.4.4 and Theorem 10.3.2.

The symbol,  $\int_a^b f(x) dx$  when  $a > b$  has not yet been defined.

**Definition 10.4.6** *Let  $[a, b]$  be an interval and let  $f \in R([a, b])$ . Then*

$$\int_b^a f(x) dx \equiv - \int_a^b f(x) dx.$$

Note that with this definition,

$$\int_a^a f(x) dx = - \int_a^a f(x) dx$$

and so

$$\int_a^a f(x) dx = 0.$$

**Theorem 10.4.7** *Assuming all the integrals make sense,*

$$\int_a^b f(x) dx + \int_b^c f(x) dx = \int_a^c f(x) dx.$$

**Proof:** This follows from Theorem 10.4.4 and Definition 10.4.6. For example, assume

$$c \in (a, b).$$

Then from Theorem 10.4.4,

$$\int_a^c f(x) dx + \int_c^b f(x) dx = \int_a^b f(x) dx$$

and so by Definition 10.4.6,

$$\begin{aligned}\int_a^c f(x) dx &= \int_a^b f(x) dx - \int_c^b f(x) dx \\ &= \int_a^b f(x) dx + \int_b^c f(x) dx.\end{aligned}$$

The other cases are similar.

The following properties of the integral have either been established or they follow quickly from what has been shown so far.

$$\text{If } f \in R([a, b]) \text{ then if } c \in [a, b], f \in R([a, c]), \quad (10.11)$$

$$\int_a^b \alpha dx = \alpha(b-a), \quad (10.12)$$

$$\int_a^b (\alpha f + \beta g)(x) dx = \alpha \int_a^b f(x) dx + \beta \int_a^b g(x) dx, \quad (10.13)$$

$$\int_a^b f(x) dx + \int_b^c f(x) dx = \int_a^c f(x) dx, \quad (10.14)$$

$$\int_a^b f(x) dx \geq 0 \text{ if } f(x) \geq 0 \text{ and } a < b, \quad (10.15)$$

$$\left| \int_a^b f(x) dx \right| \leq \int_a^b |f(x)| dx. \quad (10.16)$$

The only one of these claims which may not be completely obvious is the last one. To show this one, note that

$$|f(x)| - f(x) \geq 0, \quad |f(x)| + f(x) \geq 0.$$

Therefore, by (10.15) and (10.13), if  $a < b$ ,

$$\int_a^b |f(x)| dx \geq \int_a^b f(x) dx$$

and

$$\int_a^b |f(x)| dx \geq - \int_a^b f(x) dx.$$

Therefore,

$$\int_a^b |f(x)| dx \geq \left| \int_a^b f(x) dx \right|.$$

If  $b < a$  then the above inequality holds with  $a$  and  $b$  switched. This implies (10.16).

## 10.5 Fundamental Theorem Of Calculus

With these properties, it is easy to prove the fundamental theorem of calculus<sup>2</sup>. Let  $f \in R([a, b])$ . Then by (10.11)  $f \in R([a, x])$  for each  $x \in [a, b]$ . The first version of the fundamental theorem of calculus is a statement about the derivative of the function

$$x \rightarrow \int_a^x f(t) dt.$$

---

<sup>2</sup>This theorem is why Newton and Leibnitz are credited with inventing calculus. The integral had been around for thousands of years and the derivative was by their time well known. However the connection between these two ideas had not been fully made although Newton's predecessor, Isaac Barrow had made some progress in this direction.

**Theorem 10.5.1** Let  $f \in R([a, b])$  and let

$$F(x) \equiv \int_a^x f(t) dt.$$

Then if  $f$  is continuous at  $x \in (a, b)$ ,

$$F'(x) = f(x).$$

**Proof:** Let  $x \in (a, b)$  be a point of continuity of  $f$  and let  $h$  be small enough that  $x + h \in [a, b]$ . Then by using (10.14),

$$h^{-1}(F(x+h) - F(x)) = h^{-1} \int_x^{x+h} f(t) dt.$$

Also, using (10.12),

$$f(x) = h^{-1} \int_x^{x+h} f(x) dt.$$

Therefore, by (10.16),

$$\begin{aligned} |h^{-1}(F(x+h) - F(x)) - f(x)| &= \left| h^{-1} \int_x^{x+h} (f(t) - f(x)) dt \right| \\ &\leq \left| h^{-1} \int_x^{x+h} |f(t) - f(x)| dt \right|. \end{aligned}$$

Let  $\varepsilon > 0$  and let  $\delta > 0$  be small enough that if  $|t - x| < \delta$ , then

$$|f(t) - f(x)| < \varepsilon.$$

Therefore, if  $|h| < \delta$ , the above inequality and (10.12) shows that

$$|h^{-1}(F(x+h) - F(x)) - f(x)| \leq |h|^{-1} \varepsilon |h| = \varepsilon.$$

Since  $\varepsilon > 0$  is arbitrary, this shows

$$\lim_{h \rightarrow 0} h^{-1}(F(x+h) - F(x)) = f(x)$$

and this proves the theorem.

Note this gives existence for the initial value problem,

$$F'(x) = f(x), \quad F(a) = 0$$

whenever  $f$  is Riemann integrable and continuous.<sup>3</sup>

The next theorem is also called the fundamental theorem of calculus.

**Theorem 10.5.2** Let  $f \in R([a, b])$  and suppose there exists an antiderivative for  $f$ ,  $G$ , such that

$$G'(x) = f(x)$$

for every point of  $(a, b)$  and  $G$  is continuous on  $[a, b]$ . Then

$$\int_a^b f(x) dx = G(b) - G(a). \quad (10.17)$$

---

<sup>3</sup>Of course it was proved that if  $f$  is continuous on a closed interval,  $[a, b]$ , then  $f \in R([a, b])$  but this is a hard theorem using the difficult result about uniform continuity.

**Proof:** Let  $P = \{x_0, \dots, x_n\}$  be a partition satisfying

$$U(f, P) - L(f, P) < \varepsilon.$$

Then

$$\begin{aligned} G(b) - G(a) &= G(x_n) - G(x_0) \\ &= \sum_{i=1}^n G(x_i) - G(x_{i-1}). \end{aligned}$$

By the mean value theorem,

$$\begin{aligned} G(b) - G(a) &= \sum_{i=1}^n G'(z_i)(x_i - x_{i-1}) \\ &= \sum_{i=1}^n f(z_i) \Delta x_i \end{aligned}$$

where  $z_i$  is some point in  $[x_{i-1}, x_i]$ . It follows, since the above sum lies between the upper and lower sums, that

$$G(b) - G(a) \in [L(f, P), U(f, P)],$$

and also

$$\int_a^b f(x) dx \in [L(f, P), U(f, P)].$$

Therefore,

$$\left| G(b) - G(a) - \int_a^b f(x) dx \right| < U(f, P) - L(f, P) < \varepsilon.$$

Since  $\varepsilon > 0$  is arbitrary, (10.17) holds. This proves the theorem.

The following notation is often used in this context. Suppose  $F$  is an antiderivative of  $f$  as just described with  $F$  continuous on  $[a, b]$  and  $F' = f$  on  $(a, b)$ . Then

$$\int_a^b f(x) dx = F(b) - F(a) \equiv F(x) \Big|_a^b.$$

**Definition 10.5.3** Let  $f$  be a bounded function defined on a closed interval  $[a, b]$  and let  $P \equiv \{x_0, \dots, x_n\}$  be a partition of the interval. Suppose  $z_i \in [x_{i-1}, x_i]$  is chosen. Then the sum

$$\sum_{i=1}^n f(z_i)(x_i - x_{i-1})$$

is known as a Riemann sum. Also,

$$\|P\| \equiv \max \{|x_i - x_{i-1}| : i = 1, \dots, n\}.$$

**Proposition 10.5.4** Suppose  $f \in R([a, b])$ . Then there exists a partition,  $P \equiv \{x_0, \dots, x_n\}$  with the property that for any choice of  $z_k \in [x_{k-1}, x_k]$ ,

$$\left| \int_a^b f(x) dx - \sum_{k=1}^n f(z_k)(x_k - x_{k-1}) \right| < \varepsilon.$$

**Proof:** Choose  $P$  such that  $U(f, P) - L(f, P) < \varepsilon$  and then both  $\int_a^b f(x) dx$  and  $\sum_{k=1}^n f(z_k)(x_k - x_{k-1})$  are contained in  $[L(f, P), U(f, P)]$  and so the claimed inequality must hold. This proves the proposition.

It is significant because it gives a way of approximating the integral.

The definition of Riemann integrability given in this chapter is also called Darboux integrability and the integral defined as the unique number which lies between all upper sums and all lower sums which is given in this chapter is called the Darboux integral. The definition of the Riemann integral in terms of Riemann sums is given next.

**Definition 10.5.5** A bounded function,  $f$  defined on  $[a, b]$  is said to be Riemann integrable if there exists a number,  $I$  with the property that for every  $\varepsilon > 0$ , there exists  $\delta > 0$  such that if

$$P \equiv \{x_0, x_1, \dots, x_n\}$$

is any partition having  $\|P\| < \delta$ , and  $z_i \in [x_{i-1}, x_i]$ ,

$$\left| I - \sum_{i=1}^n f(z_i)(x_i - x_{i-1}) \right| < \varepsilon.$$

The number  $\int_a^b f(x) dx$  is defined as  $I$ .

Thus, there are two definitions of the Riemann integral. It turns out they are equivalent which is the following theorem of Darboux.

**Theorem 10.5.6** A bounded function defined on  $[a, b]$  is Riemann integrable in the sense of Definition 10.5.5 if and only if it is integrable in the sense of Darboux. Furthermore the two integrals coincide.

The proof of this theorem is left for the exercises in Problems 10 - 12. It isn't essential that you understand this theorem so if it does not interest you, leave it out. Note that it implies that given a Riemann integrable function  $f$  in either sense, it can be approximated by Riemann sums whenever  $\|P\|$  is sufficiently small. Both versions of the integral are obsolete but entirely adequate for most applications and as a point of departure for a more up to date and satisfactory integral. The reason for using the Darboux approach to the integral is that all the existence theorems are easier to prove in this context.

## 10.6 Exercises

1. Let  $F(x) = \int_{x^2}^{x^3} \frac{t^5+7}{t^7+87t^6+1} dt$ . Find  $F'(x)$ .
2. Let  $F(x) = \int_2^x \frac{1}{1+t^4} dt$ . Sketch a graph of  $F$  and explain why it looks the way it does.
3. Let  $a$  and  $b$  be positive numbers and consider the function,

$$F(x) = \int_0^{ax} \frac{1}{a^2+t^2} dt + \int_b^{a/x} \frac{1}{a^2+t^2} dt.$$

Show that  $F$  is a constant.

4. Solve the following initial value problem from ordinary differential equations which is to find a function  $y$  such that

$$y'(x) = \frac{x^7+1}{x^6+97x^5+7}, \quad y(10) = 5.$$



5. If  $F, G \in \int f(x) dx$  for all  $x \in \mathbb{R}$ , show  $F(x) = G(x) + C$  for some constant,  $C$ . Use this to give a different proof of the fundamental theorem of calculus which has for its conclusion  $\int_a^b f(t) dt = G(b) - G(a)$  where  $G'(x) = f(x)$ .
6. Suppose  $f$  is Riemann integrable on  $[a, b]$  and continuous. (In fact continuous implies Riemann integrable.) Show there exists  $c \in (a, b)$  such that

$$f(c) = \frac{1}{b-a} \int_a^b f(x) dx.$$

**Hint:** You might consider the function  $F(x) \equiv \int_a^x f(t) dt$  and use the mean value theorem for derivatives and the fundamental theorem of calculus.

7. Suppose  $f$  and  $g$  are continuous functions on  $[a, b]$  and that  $g(x) \neq 0$  on  $(a, b)$ . Show there exists  $c \in (a, b)$  such that

$$f(c) \int_a^b g(x) dx = \int_a^b f(x) g(x) dx.$$

**Hint:** Define  $F(x) \equiv \int_a^x f(t) g(t) dt$  and let  $G(x) \equiv \int_a^x g(t) dt$ . Then use the Cauchy mean value theorem on these two functions.

8. Consider the function

$$f(x) \equiv \begin{cases} \sin\left(\frac{1}{x}\right) & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases}.$$

Is  $f$  Riemann integrable? Explain why or why not.

9. Prove the second part of Theorem 10.3.2 about decreasing functions.
10. Suppose  $f$  is a bounded function defined on  $[a, b]$  and  $|f(x)| < M$  for all  $x \in [a, b]$ . Now let  $Q$  be a partition having  $n$  points,  $\{x_0^*, \dots, x_n^*\}$  and let  $P$  be any other partition. Show that

$$|U(f, P) - L(f, P)| \leq 2Mn \|P\| + |U(f, Q) - L(f, Q)|.$$

**Hint:** Write the sum for  $U(f, P) - L(f, P)$  and split this sum into two sums, the sum of terms for which  $[x_{i-1}, x_i]$  contains at least one point of  $Q$ , and terms for which  $[x_{i-1}, x_i]$  does not contain any points of  $Q$ . In the latter case,  $[x_{i-1}, x_i]$  must be contained in some interval,  $[x_{k-1}^*, x_k^*]$ . Therefore, the sum of these terms should be no larger than  $|U(f, Q) - L(f, Q)|$ .

11.  $\uparrow$  If  $\varepsilon > 0$  is given and  $f$  is a Darboux integrable function defined on  $[a, b]$ , show there exists  $\delta > 0$  such that whenever  $\|P\| < \delta$ , then

$$|U(f, P) - L(f, P)| < \varepsilon.$$

12.  $\uparrow$  Prove Theorem 10.5.6.

## 10.7 Return Of The Wild Assumption

The entire treatment of exponential functions and logarithms up till this time has been based on the Wild Assumption on Page 143. This was a totally unjustified assumption that exponential functions existed and were differentiable. Some pictures were drawn by a

computer as evidence. Based on this outrageous wild assumption, logarithms were defined. It is now possible to establish Wild Assumption 7.3.1.

Define

$$L_1(x) \equiv \int_1^x \frac{1}{t} dt. \quad (10.18)$$

There is no problem in writing this integral because the function,  $f(t) = 1/t$  is decreasing.

**Theorem 10.7.1** *The function,  $L_1 : (0, \infty) \rightarrow \mathbb{R}$  satisfies the following properties.*

$$L_1(xy) = L_1(x) + L_1(y), \quad L_1(1) = 0, \quad (10.19)$$

*The function,  $L_1$  is one to one and onto, strictly increasing, and its graph is concave downward. In addition to this, whenever  $\frac{m}{n} \in \mathbb{Q}$*

$$L_1\left(\sqrt[n]{x^m}\right) = \frac{m}{n} L_1(x). \quad (10.20)$$

**Proof:** Fix  $y > 0$  and let

$$f(x) = L_1(xy) - (L_1(x) + L_1(y))$$

Then by Theorem 6.2.6 on Page 122 and the Fundamental theorem of calculus, Theorem 10.5.1,

$$f'(x) = y \left( \frac{1}{xy} \right) - \frac{1}{x} = 0.$$

Therefore, by Corollary 6.8.4 on Page 133,  $f(x)$  is a constant. However,  $f(1) = 0$  and so this proves (10.19).

From the Fundamental theorem of calculus, Theorem 10.5.1,  $L_1'(x) = \frac{1}{x} > 0$  and so  $L_1$  is a strictly increasing function and is therefore one to one. Also the second derivative equals

$$L_1''(x) = \frac{-1}{x^2} < 0$$

showing that the graph of  $L_1$  is concave down.

Now consider the assertion that  $L_1$  is onto. First note that from the definition,  $L_1(2) > 0$ . In fact,

$$L_1(2) \geq 1/2$$

as can be seen by looking at a lower sum for  $\int_1^2 (1/t) dt$ . Now if  $x > 0$

$$L_1(x \times x) = L_1x + L_1x = 2L_1x.$$

Also,

$$\begin{aligned} L_1((x)(x)(x)) &= L_1((x)(x)) + L_1(x) \\ &= L_1(x) + L_1(x) + L_1(x) \\ &= 3L_1(x). \end{aligned}$$

Continuing in this way it follows that for any positive integer,  $n$ ,

$$L_1(x^n) = nL_1(x). \quad (10.21)$$

Therefore,  $L_1(x)$  achieves arbitrarily large values as  $x$  gets increasingly large because you can take  $x = 2$  in (10.21) and use the definition of  $L_1$  to verify that  $L_1(2) > 0$ . Now if  $x > 0$

$$0 = L_1 \left( \overbrace{\left( \frac{1}{x} \right) (x)}^{=1} \right) = L_1 \left( \frac{1}{x} \right) + L_1(x)$$

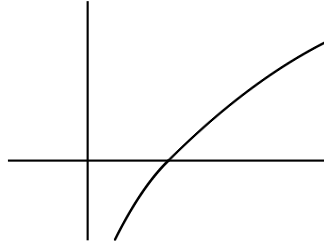
showing that

$$L_1(x^{-1}) = L_1 \left( \frac{1}{x} \right) = -L_1(x). \quad (10.22)$$

You see that  $\left(\frac{1}{2}\right)^n$  can be made as close to zero as desired by taking  $n$  sufficiently large. Also, from (10.21),

$$L_1 \left( \frac{1}{2^n} \right) = nL_1 \left( \frac{1}{2} \right) = -nL_1(2)$$

showing that  $L_1(x)$  gets arbitrarily large in the negative direction provided that  $x$  is a sufficiently small positive number. Since  $L_1$  is continuous, the intermediate value theorem may be used to fill in all the numbers in between. Thus the picture of the graph of  $L_1$  looks like the following



where the graph approaches the  $y$  axis as  $x$  gets close to 0.

It only remains to verify the claim about raising  $x$  to a rational power. From (10.21) and (10.22), for any integer,  $n$ , positive or negative,

$$L_1(x^n) = nL_1(x).$$

Therefore, letting  $m, n$  be integers,

$$mL_1(x) = L_1(x^m) = L_1 \left( \left( \sqrt[n]{x^m} \right)^n \right) = nL_1 \left( \sqrt[n]{x^m} \right)$$

and so

$$\frac{m}{n}L_1(x) = L_1 \left( \sqrt[n]{x^m} \right).$$

This proves the theorem.

Now Wild Assumption 7.3.1 on Page 143 can be fully justified.

**Definition 10.7.2** For  $b > 0$  define

$$\exp_b(x) \equiv L_1^{-1}(xL_1(b)). \quad (10.23)$$

**Proposition 10.7.3** The function just defined in (10.23) satisfies all conditions of Wild Assumption 7.3.1 on Page 143 and  $L_1(b) = \ln b$  as defined on Page 143.

**Proof:** First let  $x = \frac{m}{n}$  where  $m$  and  $n$  are integers. To verify that  $\exp_b(m/n) = \sqrt[n]{b^m}$ ,

$$\exp_b\left(\frac{m}{n}\right) \equiv L_1^{-1}\left(\frac{m}{n}L_1(b)\right) = L_1^{-1}\left(L_1\left(\sqrt[n]{b^m}\right)\right) = \sqrt[n]{b^m}$$

by (10.20).

That  $\exp_b(x) > 0$  follows immediately from the definition of the inverse function. ( $L_1$  is defined on positive real numbers and so  $L_1^{-1}$  has values in the positive real numbers.)

Consider the claim that if  $h \neq 0$  and  $b \neq 1$ , then  $\exp_b(h) \neq 1$ . Suppose then that  $\exp_b(h) = 1$ . Then doing  $L_1$  to both sides,  $hL_1(b) = L_1(1) = 0$ . Hence either  $h = 0$  or  $b = 1$  which are both excluded.

What of the laws of exponents for arbitrary values of  $x$  and  $y$ ? As part of Wild Assumption 7.3.1, these were assumed to hold.

$$L_1(\exp_b(x+y)) = (x+y)L_1(b)$$

and

$$\begin{aligned} L_1(\exp_b(x)\exp_b(y)) &= L_1(\exp_b(x)) + L_1(\exp_b(y)) \\ &= xL_1(b) + yL_1(b) = (x+y)L_1(b). \end{aligned}$$

Since  $L_1$  is one to one, this shows the first law of exponents holds,

$$\exp_b(x+y) = \exp_b(x)\exp_b(y).$$

From the definition,  $\exp_b(1) = b$  and  $\exp_b(0) = 1$ . Therefore,

$$1 = \exp_b(1 + (-1)) = \exp_b(-1)\exp_b(1) = \exp_b(-1)b$$

showing that  $\exp_b(-1) = b^{-1}$ . Now

$$L_1(\exp_{ab}(x)) = xL_1(ab) = xL_1(a) + xL_1(b)$$

while

$$\begin{aligned} L_1(\exp_a(x)\exp_b(x)) &= L_1(\exp_a(x)) + L_1(\exp_b(x)) \\ &= xL_1(a) + xL_1(b). \end{aligned}$$

Again, since  $L_1$  is one to one,  $\exp_{ab}(x) = \exp_a(x)\exp_b(x)$ . Finally,

$$L_1\left(\exp_{\exp_b(x)}(y)\right) = yL_1(\exp_b(x)) = yxL_1(b)$$

while

$$L_1(\exp_b(xy)) = xyL_1(b)$$

and so since  $L_1$  is one to one,  $\exp_{\exp_b(x)}(y) = \exp_b(xy)$ . This establishes all the laws of exponents for arbitrary real values of the exponent.

$\exp'_b$  exists by Theorem 8.1.6 on Page 151. Therefore,  $\exp_b$  defined in (10.23) satisfies all the conditions of Wild Assumption 7.3.1.

It remains to consider the derivative of  $\exp_b$  and verify  $L_1(b) = \ln b$ . First,

$$L_1(L_1^{-1}(x)) = x$$

and so

$$L'_1(L_1^{-1}(x))(L_1^{-1})'(x) = 1.$$

By (10.18),

$$\frac{(L_1^{-1})'(x)}{L_1^{-1}(x)} = 1$$

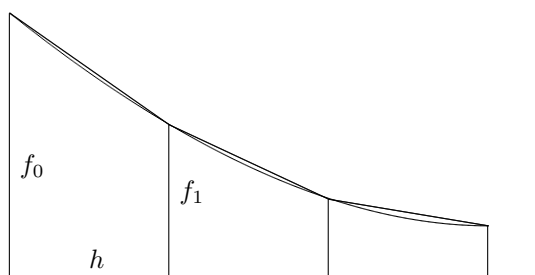
showing that  $(L_1^{-1})'(x) = (L_1^{-1})(x)$ . Now by the chain rule,

$$\begin{aligned} \exp'_b(x) &= (L_1^{-1})'(xL_1(b)) L_1(b) \\ &= L_1^{-1}(xL_1(b)) L_1(b) \\ &= \exp_b(x) L_1(b). \end{aligned}$$

From the definition of  $\ln b$  on Page 143, it follows  $\ln = L_1$ .

## 10.8 Exercises

1. Show  $\ln 2 \in [.5, 1]$ .
2. There is a general procedure for estimating the integral of a function,  $f(x) = y$  on an interval,  $[a, b]$ . Form a uniform partition,  $P = \{x_0, x_1, \dots, x_n\}$  where for each  $j$ ,  $x_j - x_{j-1} = h$ . Let  $f_i = f(x_i)$  and assuming  $f \geq 0$  on the interval  $[x_{i-1}, x_i]$ , approximate the area above this interval and under the curve with the area of a trapezoid having vertical sides,  $f_{i-1}$ , and  $f_i$  as shown in the following picture.



Thus  $\frac{1}{2} \left( \frac{f_i + f_{i-1}}{2} \right)$  approximates the area under the curve. Show that adding these up yields

$$\frac{h}{2} [f_0 + 2f_1 + \dots + 2f_{n-1} + f_n]$$

as an approximation to  $\int_a^b f(x) dx$ . This is known as the trapezoidal rule. Verify that if  $f(x) = mx + b$ , the trapezoidal rule gives the exact answer for the integral. Would this be true of upper and lower sums for such a function? Can you show that in the case of the function,  $f(t) = 1/t$  the trapezoidal rule will always yield an answer which is too large for  $\ln 2$ ?

3. Apply the trapezoidal rule to estimate  $\ln 2$  in the case where  $h = 1/5$ . Now use a calculator or table to find the exact value of  $\ln 2$ .

4. Suppose it is desired to find a function,  $L : (0, \infty) \rightarrow \mathbb{R}$  which satisfies

$$L(xy) = Lx + Ly, \quad L(1) = 0. \quad (10.24)$$

Show the only differentiable solutions to this equation are functions of the form  $L_k(x) = \int_1^x \frac{k}{t} dt$ . **Hint:** Fix  $x > 0$  and differentiate both sides of the above equation with respect to  $y$ . Then let  $y = 1$ .

5. Recall that  $\ln e = 1$ . In fact, this was how  $e$  was defined in Problem ?? on Page ??. Show that

$$\lim_{y \rightarrow 0+} (1 + yx)^{1/y} = e^x.$$

**Hint:** Consider  $\ln(1 + yx)^{1/y} = \frac{1}{y} \ln(1 + yx) = \frac{1}{y} \int_1^{1+yx} \frac{1}{t} dt$ , use upper and lower sums and then the squeezing theorem to verify  $\ln(1 + yx)^{1/y} \rightarrow x$ . Recall that  $x \rightarrow e^x$  is continuous.

6. Let there be three equally spaced points,  $x_{i-1}, x_{i-1} + h \equiv x_i$ , and  $x_i + 2h \equiv x_{i+1}$ . Suppose also a function,  $f$ , has the value  $f_{i-1}$  at  $x$ ,  $f_i$  at  $x + h$ , and  $f_{i+1}$  at  $x + 2h$ . Then consider

$$g_i(x) \equiv \frac{f_{i-1}}{2h^2} (x - x_i)(x - x_{i+1}) - \frac{f_i}{h^2} (x - x_{i-1})(x - x_{i+1}) + \frac{f_{i+1}}{2h^2} (x - x_{i-1})(x - x_i).$$

Check that this is a second degree polynomial which equals the values  $f_{i-1}, f_i$ , and  $f_{i+1}$  at the points  $x_{i-1}, x_i$ , and  $x_{i+1}$  respectively. The function,  $g_i$  is an approximation to the function,  $f$  on the interval  $[x_{i-1}, x_{i+1}]$ . Also,

$$\int_{x_{i-1}}^{x_{i+1}} g_i(x) dx$$

is an approximation to  $\int_{x_{i-1}}^{x_{i+1}} f(x) dx$ . Show  $\int_{x_{i-1}}^{x_{i+1}} g_i(x) dx$  equals

$$\frac{hf_{i-1}}{3} + \frac{hf_i 4}{3} + \frac{hf_{i+1}}{3}.$$

Now suppose  $n$  is even and  $\{x_0, x_1, \dots, x_n\}$  is a partition of the interval,  $[a, b]$  and the values of a function,  $f$  defined on this interval are  $f_i = f(x_i)$ . Adding these approximations for the integral of  $f$  on the succession of intervals,

$$[x_0, x_2], [x_2, x_4], \dots, [x_{n-2}, x_n],$$

show that an approximation to  $\int_a^b f(x) dx$  is

$$\frac{h}{3} [f_0 + 4f_1 + 2f_2 + 4f_3 + 2f_4 + \dots + 4f_{n-1} + f_n].$$

This is called Simpson's rule. Use Simpson's rule to compute an approximation to  $\ln 2$  letting  $n = 4$ . Compare with the answer from a calculator or computer.

7. Logarithms were invented before calculus, one of the inventors being Napier, a Scottish nobleman. His interest in logarithms was computational. Describe how one could use logarithms to find 7<sup>th</sup> roots for example. Also describe how one could use logarithms to do computations involving large numbers. Describe how you could construct a table for  $\log_{10} x$  for  $x$  various numbers. Next, how do you suppose Napier did it? Remember, he did not have calculus. **Hint:**  $\log_b(35^{1/7}) = \frac{1}{7} \log_b(35)$ . You can find the logarithm of the number you are after. If you had a table of logarithms to some base, you could then see what number this corresponded to.

## 10.9 Techniques Of Integration

The techniques for finding antiderivatives may be used to find integrals.

### 10.9.1 The Method Of Substitution

Recall

$$\int f(g(x)) g'(x) dx = F(x) + C, \quad (10.25)$$

where  $F'(y) = f(y)$ .

How does this relate to finding definite integrals? This is based on the following formula in which all the functions are integrable and  $F'(y) = f(y)$ .

$$\int_a^b f(g(x)) g'(x) dx = \int_{g(a)}^{g(b)} f(y) dy. \quad (10.26)$$

This formula follows from the observation that, by the fundamental theorem of calculus, both sides equal  $F(g(b)) - F(g(a))$ .

How can you remember this? The easiest way is to use the Leibniz notation. In (10.26) let  $y = g(x)$ . Then

$$\frac{dy}{dx} = g'(x)$$

and so formally  $dy = g'(x) dx$ . Then making the substitution

$$\int_a^b \overbrace{f(g(x))}^{f(y)} \overbrace{g'(x) dx}^{dy} = \int_{\text{?}}^{\text{?}} f(y) dy.$$

What should go in as the top and bottom limits of the integral? The important thing to remember is that **if you change the variable, you must change the limits!** When  $x = a$ , it follows that  $y = g(a)$  to the bottom limit must equal  $g(a)$ . Similarly the top limit should be  $g(b)$ .

**Example 10.9.1** Find  $\int_1^2 x \sin(x^2) dx$

Let  $u = x^2$  so  $du = 2x dx$  and so  $\frac{du}{2} = x dx$ . Therefore, changing the variables gives

$$\int_1^2 x \sin(x^2) dx = \frac{1}{2} \int_1^4 \sin(u) du = -\frac{1}{2} \cos(4) + \frac{1}{2} \cos(1)$$

Sometimes people prefer not to worry about the limits. This is fine provided you don't write anything which is false. The above problem can be done in the following way.

$$\int x \sin(x^2) dx = -\frac{1}{2} \cos(x^2) + C$$

and so from an application of the fundamental theorem of calculus

$$\begin{aligned} \int_1^2 x \sin(x^2) dx &= -\frac{1}{2} \cos(x^2) \Big|_1^2 \\ &= -\frac{1}{2} \cos(4) + \frac{1}{2} \cos(1). \end{aligned}$$

**Example 10.9.2** Find the area of the ellipse,

$$\frac{(y - \beta)^2}{b^2} + \frac{(x - \alpha)^2}{a^2} = 1.$$

If you sketch the ellipse, you see that it suffices to find the area of the top right quarter for  $y \geq \beta$  and  $x \geq \alpha$  and multiply by 4 since the bottom half is just a reflection of the top half about the line,  $y = \beta$  and the left top quarter is just the reflection of the top right quarter reflected about the line,  $x = \alpha$ . Thus the area of the ellipse is

$$4 \int_{\alpha}^{\alpha+a} b \sqrt{1 - \frac{(x - \alpha)^2}{a^2}} dx$$

Change the variables, letting  $u = \frac{x - \alpha}{a}$ . Then  $du = \frac{1}{a} dx$  and so upon changing the limits to correspond to the new variables, this equals

$$4ba \int_0^1 \sqrt{1 - u^2} du = 4 \times ab \times \frac{\pi/4}{1} = \pi ab$$

because the integral in the above is just one quarter of the unit circle and so has area equal to  $\pi/4$ .

## 10.9.2 Integration By Parts

Recall the following proposition for finding antiderivative.

**Proposition 10.9.3** Let  $u$  and  $v$  be differentiable functions for which  $\int u(x) v'(x) dx$  and  $\int u'(x) v(x) dx$  are nonempty. Then

$$uv - \int u'(x) v(x) dx = \int u(x) v'(x) dx. \quad (10.27)$$

In terms of integrals, this is stated in the following proposition.

**Proposition 10.9.4** Let  $u$  and  $v$  be differentiable functions on  $[a, b]$  such that  $uv', u'v \in R([a, b])$ . Then

$$\int_a^b u(x) v'(x) dx = uv(x) \Big|_a^b - \int_a^b u'(x) v(x) dx \quad (10.28)$$

**Proof:** Use the product rule and properties of integrals to write

$$\begin{aligned} \int_a^b u(x) v'(x) dx &= \int_a^b (uv)'(x) dx - \int_a^b u'(x) v(x) dx \\ &= uv(x) \Big|_a^b - \int_a^b u'(x) v(x) dx. \end{aligned}$$

This proves the proposition.

**Example 10.9.5** Find  $\int_0^\pi x \sin(x) dx$

Let  $u(x) = x$  and  $v'(x) = \sin(x)$ . Then applying (10.27),

$$\begin{aligned} \int_0^\pi x \sin(x) dx &= (-\cos(x)) x \Big|_0^\pi - \int_0^\pi (-\cos(x)) dx \\ &= -\pi \cos(\pi) = \pi. \end{aligned}$$



**Example 10.9.6** Find  $\int_0^1 x e^{2x} dx$

Let  $u(x) = x$  and  $v'(x) = e^{2x}$ . Then from (10.28)

$$\begin{aligned}\int_0^1 x e^{2x} dx &= \frac{e^{2x}}{2} x \Big|_0^1 - \int_0^1 \frac{e^{2x}}{2} dx \\ &= \frac{e^{2x}}{2} x \Big|_0^1 - \frac{e^{2x}}{4} \Big|_0^1 \\ &= \frac{e^2}{2} - \left( \frac{e^2}{4} - \frac{1}{4} \right) \\ &= \frac{e^2}{4} + \frac{1}{4}\end{aligned}$$

## 10.10 Exercises

1. Find the integrals.

(a)  $\int_0^4 x e^{-3x} dx$

(b)  $\int_2^3 \frac{1}{x(\ln(|x|))^2} dx$

(c)  $\int_0^1 x \sqrt{2-x} dx$

(d)  $\int_2^3 (\ln |x|)^2 dx$  **Hint:** Let  $u(x) = (\ln |x|)^2$  and  $v'(x) = 1$ .

(e)  $\int_0^\pi x^3 \cos(x^2) dx$

2. Find  $\int_1^2 x \ln(x^2) dx$

3. Find  $\int_0^1 e^x \sin(x) dx$

4. Find  $\int_0^1 2^x \cos(x) dx$

5. Find  $\int_0^2 x^3 \cos(x) dx$

6. Find the integrals.

(a)  $\int_5^6 \frac{x}{\sqrt{2x-3}} dx$

(b)  $\int_2^4 x (3x^2 + 6)^5 dx$

(c)  $\int_0^\pi x \sin(x^2) dx$

(d)  $\int_0^{\pi/4} \sin^3(2x) \cos(2x) dx$

(e)  $\int_0^7 \frac{1}{\sqrt{1+4x^2}} dx$  **Hint:** Remember the  $\sinh^{-1}$  function and its derivative.

7. Find the integrals.

(a)  $\int_0^{\pi/9} \sec(3x) dx$

(b)  $\int_0^{\pi/9} \sec^2(3x) \tan(3x) dx$

(c)  $\int_0^5 \frac{1}{3+5x^2} dx$

(d)  $\int_0^1 \frac{1}{\sqrt{5-4x^2}} dx$

(e)  $\int_2^6 \frac{3}{x\sqrt{4x^2-5}} dx$

8. Find the integrals.

(a)  $\int_0^3 x \cosh(x^2 + 1) dx$

(b)  $\int_0^2 x^3 5^{x^4} dx$

(c)  $\int_{-\pi}^{\pi} \sin(x) 7^{\cos(x)} dx$

(d)  $\int_0^{\pi} x \sin(x^2) dx$

(e)  $\int_1^2 x^5 \sqrt{2x^2 + 1} dx$  **Hint:** Let  $u = 2x^2 + 1$ .

9. Find  $\int_0^{\pi/4} \sin^2(x) dx$ . **Hint:** Derive and use  $\sin^2(x) = \frac{1 - \cos(2x)}{2}$ .

10. Find the area between the graphs of  $y = \sin(2x)$  and  $y = \cos(2x)$  for  $x \in [0, 2\pi]$ .

11. Find the following integrals.

(a)  $\int_1^{\pi} x^2 \sin(x^3) dx$

(b)  $\int_1^6 \frac{x}{1+x^2} dx$

(c)  $\int_0^{.5} \frac{1}{1+4x^2} dx$

(d)  $\int_1^4 x^2 \sqrt{1+x} dx$

12. The most important of all differential equations is the first order linear equation,  $y' + p(t)y = f(t)$ . Show the solution to the initial value problem consisting of this equation and the initial condition,  $y(a) = y_a$  is

$$y(t) = e^{-P(t)} y_a + e^{-P(t)} \int_a^t e^{P(s)} f(s) ds,$$

where  $P(t) = \int_a^t p(s) ds$ . Give conditions under which everything is correct. **Hint:** You use the integrating factor approach. Multiply both sides by  $e^{P(t)}$ , verify the left side equals

$$\frac{d}{dt} \left( e^{P(t)} y(t) \right),$$

and then take the integral,  $\int_a^t$  of both sides.

13. Suppose  $x_0 \in (a, b)$  and that  $f$  is a function which has  $n+1$  continuous derivatives on this interval. Consider the following.

$$\begin{aligned} f(x) &= f(x_0) + \int_{x_0}^x f'(t) dt \\ &= f(x_0) + (t-x)f'(t) \Big|_{x_0}^x + \int_{x_0}^x (x-t)f''(t) dt \\ &= f(x_0) + f'(x_0)(x-x_0) + \int_{x_0}^x (x-t)f''(t) dt. \end{aligned}$$

Explain the above steps and continue the process to eventually obtain Taylor's formula,

$$f(x) = f(x_0) + \sum_{k=1}^n \frac{f^{(k)}(x_0)}{k!} (x-x_0)^k + \frac{1}{n!} \int_{x_0}^x (x-t)^n f^{(n+1)}(t) dt$$

where  $n! \equiv n(n-1) \cdots 3 \cdot 2 \cdot 1$  if  $n \geq 1$  and  $0! \equiv 1$ .

14. In the above Taylor's formula, use Problem 7 on Page 241 to obtain the existence of some  $z$  between  $x_0$  and  $x$  such that

$$f(x) = f(x_0) + \sum_{k=1}^n \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k + \frac{f^{(n+1)}(z)}{(n+1)!} (x - x_0)^{n+1}.$$

**Hint:** You might consider two cases, the case when  $x > x_0$  and the case when  $x < x_0$ .

15. There is a general procedure for constructing these methods of approximate integration like the trapezoidal rule and Simpson's rule. Consider  $[0, 1]$  and divide this interval into  $n$  pieces using a uniform partition,  $\{x_0, \dots, x_n\}$  where  $x_i - x_{i-1} = 1/n$  for each  $i$ . The approximate integration scheme for a function,  $f$ , will be of the form

$$\left(\frac{1}{n}\right) \sum_{i=0}^n c_i f_i \approx \int_0^1 f(x) dx$$

where  $f_i = f(x_i)$  and the constants,  $c_i$  are chosen in such a way that the above sum gives the exact answer for  $\int_0^1 f(x) dx$  where  $f(x) = 1, x, x^2, \dots, x^n$ . When this has been done, change variables to write

$$\begin{aligned} \int_a^b f(y) dy &= (b-a) \int_0^1 f(a + (b-a)x) dx \\ &\approx \frac{b-a}{n} \sum_{i=1}^n c_i f\left(a + (b-a)\left(\frac{i}{n}\right)\right) \\ &= \frac{b-a}{n} \sum_{i=1}^n c_i f_i \end{aligned}$$

where  $f_i = f\left(a + (b-a)\left(\frac{i}{n}\right)\right)$ . Consider the case where  $n = 1$ . It is necessary to find constants  $c_0$  and  $c_1$  such that

$$\begin{aligned} c_0 + c_1 &= 1 = \int_0^1 1 dx \\ 0c_0 + c_1 &= 1/2 = \int_0^1 x dx. \end{aligned}$$

Show that  $c_0 = c_1 = 1/2$ , and that this yields the trapezoidal rule. Next take  $n = 2$  and show the above procedure yields Simpson's rule. Show also that if this integration scheme is applied to any polynomial of degree 3 the result will be exact. That is,

$$\frac{1}{2} \left( \frac{1}{3} f_0 + \frac{4}{3} f_1 + \frac{1}{3} f_2 \right) = \int_0^1 f(x) dx$$

whenever  $f(x)$  is a polynomial of degree three. Show that if  $f_i$  are the values of  $f$  at  $a$ ,  $\frac{a+b}{2}$ , and  $b$  with  $f_1 = f\left(\frac{a+b}{2}\right)$ , it follows that the above formula gives  $\int_a^b f(x) dx$  exactly whenever  $f$  is a polynomial of degree three. Obtain an integration scheme for  $n = 3$ .

16. Let  $f$  have four continuous derivatives on  $[x_{i-1}, x_{i+1}]$  where  $x_{i+1} = x_{i-1} + 2h$  and  $x_i = x_{i-1} + h$ . Show using Problem 14, there exists a polynomial of degree three,  $p_3(x)$ , such that

$$f(x) = p_3(x) + \frac{1}{4!} f^{(4)}(\xi) (x - x_i)^4$$

Now use Problem 15 and Problem 6 on Page 246 to conclude

$$\left| \int_{x_{i-1}}^{x_{i+1}} f(x) dx - \left( \frac{hf_{i-1}}{3} + \frac{hf_i}{3} + \frac{hf_{i+1}}{3} \right) \right| < \frac{M}{4!} \frac{2h^5}{5},$$

where  $M$  satisfies,  $M \geq \max \{|f^{(4)}(t)| : t \in [x_{i-1}, x_i]\}$ . Now let  $S(a, b, f, 2m)$  denote the approximation to  $\int_a^b f(x) dx$  obtained from Simpson's rule using  $2m$  equally spaced points. Show

$$\left| \int_a^b f(x) dx - S(a, b, f, 2m) \right| < \frac{M}{1920} (b-a)^5 \frac{1}{m^4}$$

where  $M \geq \max \{|f^{(4)}(t)| : t \in [a, b]\}$ . Better estimates are available in numerical analysis books. However, these also have the error in the form  $C(1/m^4)$ .

## 10.11 Improper Integrals

The integral is only defined for certain bounded functions which are defined on closed and bounded intervals. Nevertheless people do consider things like the following:  $\int_0^\infty f(t) dt$ . Whenever things like this occur they require a special definition. In this section a few types of improper integrals will be discussed.

**Definition 10.11.1** The symbol  $\int_a^\infty f(t) dt$  is defined to equal

$$\lim_{R \rightarrow \infty} \int_a^R f(t) dt$$

whenever this limit exists. If  $\lim_{x \rightarrow a+} f(t) = \pm\infty$  but  $f$  is integrable on  $[a + \delta, b]$  for all small  $\delta$ , then

$$\int_a^b f(t) dt \equiv \lim_{\delta \rightarrow 0+} \int_{a+\delta}^b f(t) dt$$

whenever this limit exists. Similarly, if  $\lim_{x \rightarrow b-} f(t) = \pm\infty$  but  $f$  is integrable on  $[a, b - \delta]$  for all small  $\delta$ , then

$$\int_a^b f(t) dt \equiv \lim_{\delta \rightarrow 0+} \int_a^{b-\delta} f(t) dt$$

whenever the limit exists. Finally, if  $\lim_{x \rightarrow a+} f(t) = \pm\infty$ , then

$$\int_a^\infty f(t) dt \equiv \lim_{R \rightarrow \infty} \int_a^R f(t) dt$$

where the improper integral,  $\int_a^R f(t) dt$  is defined above as  $\lim_{\delta \rightarrow 0+} \int_{a+\delta}^R f(t) dt$ .

You can probably construct other examples of improper integrals such as integrals of the form  $\int_{-\infty}^a f(t) dt$ . The definitions are analogous to the above.

**Example 10.11.2** Find  $\int_0^\infty e^{-t} dt$ .

From the definition, this equals  $\lim_{R \rightarrow \infty} \int_0^R e^{-t} dt = \lim_{R \rightarrow \infty} (1 - e^{-R}) = 1$ .

**Example 10.11.3** Find  $\int_0^1 \frac{1}{\sqrt{x}} dx$ .

From the definition this equals  $\lim_{\delta \rightarrow 0+} \int_{\delta}^1 x^{-1/2} dx = \lim_{\delta \rightarrow 0+} (2 - 2\sqrt{\delta}) = 2$ .

Sometimes you can argue the improper integral exists even though you can't find it. The following theorem is about this question of existence.

**Theorem 10.11.4** *Suppose  $f(t) \geq 0$  for all  $t \in [a, \infty)$  and that  $f \in R([a + \delta, R])$  whenever  $\delta > 0$  is small enough, for every  $R > a$ . Suppose also there exists a number,  $M$  such that for all  $R > a$ , the integral  $\int_a^R f(t) dt$  exists and*

$$\int_a^R f(t) dt \leq M. \quad (10.29)$$

*Then  $\int_a^\infty f(t) dt$  exists.*

*If  $f(t) \geq 0$  for all  $t \in (a, b]$  and there exists  $M$  such that*

$$\int_{a+\delta}^b f(t) dt \leq M \quad (10.30)$$

*for all  $\delta > 0$ , then  $\int_a^b f(t) dt$  exists.*

*If  $f(t) \geq 0$  for all  $t \in [a, b)$  and there exists  $M$  such that*

$$\int_a^{b-\delta} f(t) dt \leq M \quad (10.31)$$

*for all  $\delta > 0$ , then  $\int_a^b f(t) dt$  exists.*

**Proof:** Suppose (10.29). Then  $I \equiv \sup \left\{ \int_a^R f(t) dt : R > a \right\} \leq M$ . It follows that if  $\varepsilon > 0$  is given, there exists  $R_0$  such that  $\int_a^{R_0} f(t) dt \in (I - \varepsilon, I]$ . Then, since  $f(t) \geq 0$ , it follows that for  $R \geq R_0$ , and small  $\delta > 0$ ,

$$\int_{a+\delta}^R f(t) dt = \int_{a+\delta}^{R_0} f(t) dt + \int_{R_0}^R f(t) dt.$$

Letting  $\delta \rightarrow 0+$ ,

$$\int_a^R f(t) dt = \int_a^{R_0} f(t) dt + \int_{R_0}^R f(t) dt \geq \int_a^{R_0} f(t) dt.$$

Therefore, whenever  $R > R_0$ ,  $\left| \int_a^R f(t) dt - I \right| < \varepsilon$ . Since  $\varepsilon$  is arbitrary, the conditions for

$$I = \lim_{R \rightarrow \infty} \int_a^R f(t) dt$$

are satisfied and so  $I = \int_a^\infty f(t) dt$ .

Now suppose (10.30). Then  $I \equiv \sup \left\{ \int_{a+\delta}^b f(t) dt : \delta > 0 \right\} \leq M$ . It follows that if  $\varepsilon > 0$  is given, there exists  $\delta_0 > 0$  such that

$$\int_{a+\delta_0}^b f(t) dt \in (I - \varepsilon, I].$$

Therefore, if  $\delta < \delta_0$ ,

$$I - \varepsilon < \int_{a+\delta_0}^b f(t) dt \leq \int_{a+\delta}^b f(t) dt \leq I$$

showing that for such  $\delta$ ,  $\left| \int_{a+\delta}^b f(t) dt - I \right| < \varepsilon$ . This is what is meant by the expression

$$\lim_{\delta \rightarrow 0+} \int_{a+\delta}^b f(t) dt = I$$

and so  $I = \int_a^b f(t) dt$ .

The last case is entirely similar to this one. This proves the theorem.

**Example 10.11.5** Does  $\int_0^1 \frac{1}{\sqrt{\sin x}} dx$  exist?

I don't know how to find an antiderivative for this function but the question of existence can still be resolved. Since  $\lim_{x \rightarrow 0+} \frac{x}{\sin x} = 1$ , it follows that for  $x$  small enough,  $\frac{x}{\sin x} < \frac{3}{2}$ , say for  $x < \delta_1$ . Then for such  $x$ , it follows

$$\frac{2}{3}x < \sin x$$

and so if  $\delta < \delta_1$ ,

$$\int_{\delta}^1 \frac{1}{\sqrt{\sin x}} dx \leq \int_{\delta}^{\delta_1} \sqrt{\frac{3}{2}} \frac{1}{\sqrt{x}} dx + \int_{\delta_1}^1 \frac{1}{\sqrt{\sin \delta_1}} dx.$$

Now using the argument of Example 10.11.3, the first integral in the above is bounded above by  $(\sqrt{\delta_1} - \sqrt{\delta})\sqrt{6}$ . The second integral equals  $\frac{1-\delta_1}{\sqrt{\sin \delta_1}}$ . Therefore, the improper integral exists because the conditions of Theorem 10.11.4 with  $M = \frac{1-\delta_1}{\sqrt{\sin \delta_1}} + (\sqrt{\delta_1} - \sqrt{\delta})\sqrt{6}$ .

**Example 10.11.6** The gamma function is defined by  $\Gamma(\alpha) \equiv \int_0^{\infty} e^{-t} t^{\alpha-1} dt$  whenever  $\alpha > 0$ . Does the improper integral exist?

You should supply the details to the following estimate in which  $\delta$  is a small positive number less than 1 and  $R$  is a large positive number.

$$\begin{aligned} \int_{\delta}^R e^{-t} t^{\alpha-1} dt &\leq \int_{\delta}^k e^{-t} t^{\alpha-1} dt + \int_k^R e^{-t} t^{\alpha-1} dt \\ &\leq \int_0^k t^{\alpha-1} dt + \int_k^{\infty} e^{-t/2} dt. \end{aligned}$$

Here  $k$  is chosen such that if  $t \geq k$ ,

$$e^{-t} t^{\alpha-1} < e^{-t/2}.$$

Such a  $k$  exists because

$$\lim_{t \rightarrow \infty} \frac{e^{-t} t^{\alpha-1}}{e^{-t/2}} = 0.$$

Therefore, let  $M \equiv \int_0^k t^{\alpha-1} dt + \int_k^{\infty} e^{-t/2} dt$  and this shows from the above theorem that

$$\int_0^R e^{-t} t^{\alpha-1} dt \leq M$$

for all large  $R$  and so  $\int_0^{\infty} e^{-t} t^{\alpha-1} dt$  exists.

Sometimes the existence of the improper integral is a little more subtle. This is the case when functions are not all the same sign for example.

**Example 10.11.7** Does  $\int_0^\infty \frac{\sin x}{x} dx$  exist?

You should verify  $\int_0^1 \frac{\sin x}{x} dx$  exists and that

$$\begin{aligned}\int_0^R \frac{\sin x}{x} dx &= \int_0^1 \frac{\sin x}{x} dx + \int_1^R \frac{\sin x}{x} dx \\ &= \int_0^1 \frac{\sin x}{x} dx + \cos 1 - \frac{\cos R}{R} - \int_1^R \frac{\cos x}{x^2} dx.\end{aligned}$$

Thus the improper integral exists if it can be shown that  $\int_1^\infty \frac{\cos x}{x^2} dx$  exists. However,

$$\int_1^R \frac{\cos x}{x^2} dx = \int_1^R \frac{\cos x + |\cos x|}{x^2} dx - \int_1^R \frac{(|\cos x| - \cos x)}{x^2} dx \quad (10.32)$$

and

$$\begin{aligned}\int_1^R \frac{\cos x + |\cos x|}{x^2} dx &\leq \int_1^R \frac{2}{x^2} dx \leq \int_0^\infty \frac{2}{x^2} dx < \infty \\ \int_1^R \frac{(|\cos x| - \cos x)}{x^2} dx &\leq \int_1^R \frac{2}{x^2} dx \leq \int_0^\infty \frac{2}{x^2} dx < \infty.\end{aligned}$$

Since both integrands are positive, Theorem 10.11.4 applies and the limits

$$\lim_{R \rightarrow \infty} \int_1^R \frac{(|\cos x| - \cos x)}{x^2} dx, \quad \lim_{R \rightarrow \infty} \int_1^R \frac{\cos x + |\cos x|}{x^2} dx$$

both exist and so from (10.32)  $\lim_{R \rightarrow \infty} \int_1^R \frac{\cos x}{x^2} dx$  also exists and so  $\int_0^\infty \frac{\sin x}{x} dx$  exists.

This is an important example. There are at least two ways to show that  $\int_0^\infty \frac{\sin x}{x} dx = \frac{1}{2}\pi$ . However, they involve techniques which will not be discussed in this book. It is a standard problem in the subject of complex analysis. The above argument is a special case of the following corollary to Theorem 10.11.4.

**Definition 10.11.8** Let  $f$  be a real valued function. Then  $f^+(t) \equiv \frac{|f(t)| + f(t)}{2}$  and  $f^-(t) \equiv \frac{|f(t)| - f(t)}{2}$ . Thus  $|f(t)| = f^+(t) + f^-(t)$  and  $f(t) = f^+(t) - f^-(t)$  while both  $f^+$  and  $f^-$  are nonnegative functions.

**Corollary 10.11.9** Suppose  $f$  is a real valued function and the conditions of Theorem 10.11.4 hold for both  $f^+$  and  $f^-$ . Then  $\int_0^\infty f(t) dt$  exists.

**Example 10.11.10** Does  $\int_0^\infty \cos(x^2) dx$  exist?

This is called a Fresnel integral and it has also been evaluated exactly using techniques from complex analysis. In fact  $\int_0^\infty \cos(x^2) dx = \frac{1}{4}\sqrt{2}\sqrt{\pi}$ . The verification that this integral exists is left to you. First change the variable letting  $x^2 = u$  and then integrate by parts. You will eventually get an integral of the form  $\int_0^\infty \frac{\sin u}{u^{3/2}} du$ . Break it up into positive and negative parts and use Corollary 10.11.9.

## 10.12 Exercises

1. Verify all the details in Example 10.11.6.
2. Verify all the details of Example 10.11.7.

3. Verify all the details of Example 10.11.10.
4. Show  $\Gamma(1) = 1 = \Gamma(2)$ . Next show  $\Gamma(\alpha + 1) = \alpha\Gamma(\alpha)$  and prove that for  $n$  a non-negative integer,  $\Gamma(n + 1) = n!$ .
5. It can be shown that  $\Gamma(\frac{1}{2}) = \sqrt{\pi}$ . Using this and Problem 4, find  $\Gamma(\frac{5}{2})$ . What is the advantage of the gamma function over the notion of factorials? **Hint:** For what values of  $x$  is  $\Gamma(x)$  defined?
6. Prove  $\int_0^\infty \sin(x^2) dx$  exists.
7. For  $\alpha > 0$  find  $\int_0^1 t^\alpha \ln t dt$  if it exists and if it does not exist, explain why.
8. Prove that for every  $\alpha > 0$ ,  $\int_0^1 t^{\alpha-1} dt$  exists and find the answer.
9. Prove that for every  $\alpha > 0$ ,  $\int_0^1 (\tan t)^{\alpha-1} dt$  exists.
10. Prove that for every  $\alpha > 0$ ,  $\int_0^1 (\sin t)^{\alpha-1} dt$  exists.
11. Recall the area of the surface obtained by revolving the graph of  $y = f(x)$  about the  $x$  axis for  $x \in [a, b]$  for  $f$  a positive continuous function having continuous derivative is given by

$$2\pi \int_a^b f(x) \sqrt{1 + (f'(x))^2} dx.$$

Also the volume of the solid obtained in the same way is given by the integral,

$$\pi \int_a^b (f(x))^2 dx.$$

It seems reasonable to define the surface area and volume for  $x \in [a, \infty)$  in terms of an improper integral. Try it on the function,  $f(x) = 1/x$  for  $x \geq 1$ . Show the resulting solid has finite volume but infinite surface area. Thus you could fill it but you couldn't paint it. Sometimes people call this Gabriel's horn.

12. Show  $\int_0^\infty \frac{2x}{1+x^2} dx$  does not exist but that  $\lim_{R \rightarrow \infty} \int_{-R}^R \frac{2x}{1+x^2} dx = 0$ . This last limit is called the Cauchy principle value integral. It is not a very respectable thing. The methods of complex analysis typically give the Cauchy principle value. Try to show the following: For every number,  $A$ , there exist sequences  $a_n, b_n \rightarrow \infty$  such that

$$\lim_{n \rightarrow \infty} \int_{-a_n}^{b_n} \frac{2x}{1+x^2} dx = A.$$

This is true and shows why such principle value integrals are somewhat disreputable.

13. Suppose  $f$  is a continuous function which is bounded and defined on  $\mathbb{R}$ . Show

$$\int_{-\infty}^{\infty} \frac{\varepsilon}{\pi(\varepsilon^2 + (x - x_1)^2)} dx = 1.$$

Next show that

$$\lim_{\varepsilon \rightarrow 0+} \int_{-\infty}^{\infty} \frac{\varepsilon f(x)}{\pi(\varepsilon^2 + (x - x_1)^2)} dx = f(x_1).$$



# Infinite Series

## 11.1 Approximation By Taylor Polynomials

By now, you have noticed there are two sorts of functions, those which come from a formula like  $f(x) = x^2 + 2$  which are easy to evaluate by following a simple procedure, and those which come as short words; things like  $\ln(x)$  or  $\sin(x)$ . This latter type of function is not so easy to evaluate. For example, what is  $\sin 2$ ? Can you get it by doing a simple sequence of operations like you can with  $f(x) = x^2 + 2$ ? How can you find  $\sin 2$ ? It turns out there are many ways to do so. In this section, the method of Taylor polynomials is discussed. The following theorem is called Taylor's theorem. Before presenting it, recall the meaning of  $n!$  for  $n$  a positive integer. Define  $0! \equiv 1 = 1!$  and  $(n+1)! \equiv (n+1)n!$  so that  $n! = n(n-1) \cdots 1$ . In particular,  $2! = 2$ ,  $3! = 3 \times 2! = 6$ ,  $4! = 4 \times 3! = 24$ , etc.

**Theorem 11.1.1** *Suppose  $f$  has  $n+1$  derivatives on an interval,  $(a, b)$  and let  $c \in (a, b)$ . Then if  $x \in (a, b)$ , there exists  $\xi$  between  $c$  and  $x$  such that*

$$f(x) = f(c) + \sum_{k=1}^n \frac{f^{(k)}(c)}{k!} (x-c)^k + \frac{f^{(n+1)}(\xi)}{(n+1)!}.$$

(In this formula, the symbol  $\sum_{k=1}^0 a_k$  will denote the number 0.)

**Proof:** Suppose first that  $n = 0$ . By the mean value theorem,

$$f(x) - f(c) = f'(\xi)(x-c)$$

for some  $\xi$  between  $c$  and  $x$ . Now assume the theorem holds for  $n$ . Then let

$$g(x) \equiv f(x) - f(c) - \sum_{k=1}^{n+1} \frac{f^{(k)}(c)}{k!} (x-c)^k, \quad h(x) = (x-c)^{n+2}$$

so  $g(c) = h(c) = 0$ . By the Cauchy mean value theorem, Theorem 6.8.2 on Page 133, and the observation that  $g(c)$  and  $h(c)$  both equal 0, there exists  $\xi$  between  $x$  and  $c$  such that  $h'(\xi)g(x) = g'(\xi)h(x)$  and so

$$\begin{aligned} (n+2)(\xi-c)^{n+1} \left( f(x) - f(c) - \sum_{k=1}^{n+1} \frac{f^{(k)}(c)}{k!} (x-c)^k \right) = \\ \left( f'(\xi) - \sum_{k=1}^{n+1} \frac{f^{(k)}(c)}{(k-1)!} (\xi-c)^{k-1} \right) (x-c)^{n+2}. \end{aligned} \quad (11.1)$$

The induction hypothesis in (11.1) yields

$$\begin{aligned} f'(\xi) - \sum_{k=1}^{n+1} \frac{f^{(k)}(c)}{(k-1)!} (\xi - c)^{k-1} &= f'(\xi) - f'(c) - \sum_{k=1}^n \frac{(f')^{(k)}(c)}{k!} (\xi - c)^k \\ &= \frac{(f')^{(n+1)}(\xi_1)}{(n+1)!} (\xi - c)^{n+1} \end{aligned}$$

for some  $\xi_1$  between  $\xi$  and  $c$ . Consequently,

$$\begin{aligned} (n+2)(\xi - c)^{n+1} \left( f(x) - f(c) - \sum_{k=1}^{n+1} \frac{f^{(k)}(c)}{k!} (x - c)^k \right) &= \\ \left( \frac{(f')^{(n+1)}(\xi_1)}{(n+1)!} (\xi - c)^{n+1} \right) (x - c)^{n+2}. \end{aligned} \quad (11.2)$$

and so, dividing by  $(n+2)(\xi - c)^{n+1}$ ,

$$f(x) - f(c) - \sum_{k=1}^{n+1} \frac{f^{(k)}(c)}{k!} = \frac{(f')^{(n+1)}(\xi_1)}{(n+2)!} (x - c)^{n+2}$$

where  $\xi_1$  is between  $c$  and  $\xi$  while  $\xi$  is between  $c$  and  $x$ . Thus  $\xi_1$  is between  $c$  and  $x$  and this proves the theorem.

The term  $\frac{f^{(n+1)}(\xi)}{(n+1)!}$ , is called the remainder and this particular form of the remainder is called the Lagrange form of the remainder.

**Example 11.1.2** Approximate  $\sin x$  for  $x$  in some open interval containing 0.

Use Taylor's formula just presented and let  $c = 0$ . Then for  $f(x) = \sin x$ ,

$$f'(x) = \cos x, \quad f''(x) = -\sin x, \quad f'''(x) = -\cos x,$$

etc. Therefore,  $f(0) = 0, f'(0) = 1, f''(0) = 0, f'''(0) = -1$ , etc. Thus the Taylor polynomial for  $\sin x$  is of the form

$$x - \frac{x^3}{3!} + \cdots \pm \frac{x^{2n+1}}{(2n+1)!}$$

while the remainder is of the form

$$\frac{f^{(2n+2)}(\xi)}{(2n+2)!}$$

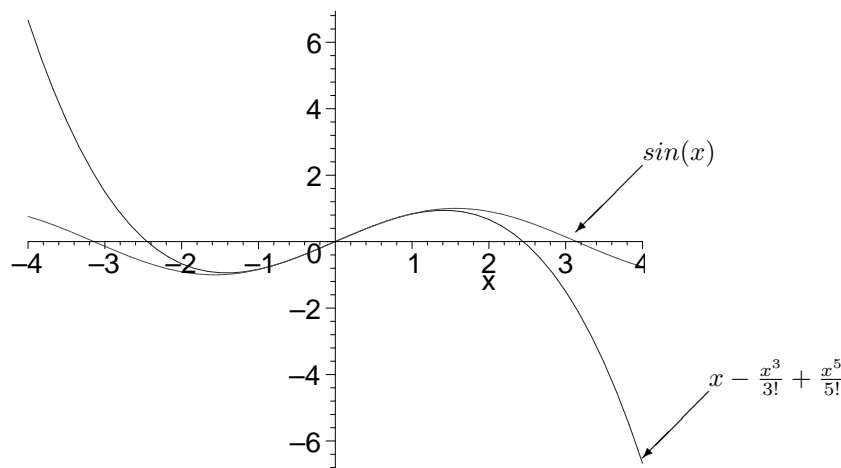
for some  $\xi$  between 0 and  $x$ . For  $n = 2$  in the above, the resulting polynomial is

$$x - \frac{x^3}{3!} + \frac{x^5}{5!}$$

and the error between this polynomial and  $\sin x$  must be measured by the remainder term. Therefore,

$$\left| \sin x - \left( x - \frac{x^3}{3!} + \frac{x^5}{5!} \right) \right| \leq \left| \frac{f^{(6)}(\xi) x^6}{6!} \right| \leq \frac{x^6}{6!}.$$

For small  $x$ , this error is very small but if  $x$  is large, no such conclusion can be drawn. This is illustrated in the following picture.



You see from the picture that the polynomial is a very good approximation for the function,  $\sin x$  as long as  $|x|$  is small but that if  $|x|$  gets very large, the approximation is lousy. The above estimate indicates the good approximation holds as long as  $|x|$  is small and it quantifies how good the approximation is. Suppose for example, you wanted to find  $\sin(.5)$ . Then from the above error estimate,

$$\left| \sin(.5) - \left( (.5) - \frac{(.5)^3}{3!} + \frac{(.5)^5}{5!} \right) \right| \leq \frac{(.5)^6}{6!} = \frac{1}{46080}$$

so difference between the approximation and  $\sin(.5)$  is less than  $10^{-4}$ . If this is used to find  $\sin(.1)$  the polynomial approximation would be even closer.

## 11.2 Exercises

1. Let  $p_n(x) = a_0 + \sum_{k=1}^n a_k(x-c)^k$ . Show that if you require that  $p_n(c) = f(c)$ ,  $p'_n(c) = f'(c)$ ,  $\dots$ ,  $p_n^{(n)}(c) = f^{(n)}(c)$ , then this requirement is achieved if and only if  $a_0 = f(c)$ ,  $a_1 = f'(c)$ ,  $\dots$ ,  $a_n = \frac{f^{(n)}(c)}{n!}$ . Thus the Taylor polynomial of degree  $n$  and its first  $n$  derivatives agree with the function and its first  $n$  derivatives when  $x = c$ .
2. Find the Taylor polynomials for  $\cos x$  for  $x$  near 0 along with a formula for the remainder. Use your approximate polynomial to compute  $\cos(.5)$  to 3 decimal places and prove your approximation is this good.
3. Find the Taylor polynomials for  $\ln(1+x)$  for  $x$  near 0.
4. Find the Taylor polynomials for  $\ln(1-x)$  for  $x$  near 0.
5. Find a Taylor polynomial for  $\ln\left(\frac{1+x}{1-x}\right)$  and use it to compute  $\ln 5$  to three decimal places.

6. Verify that  $\lim_{n \rightarrow \infty} \frac{M^n}{n!} = 0$  whenever  $M$  is a positive real number. **Hint:** Prove by induction that  $M^n/n! \leq (2M)^n/(\sqrt{n})^n$ . Now consider what happens when  $\sqrt{n}$  is much larger than  $2M$ .
7. Show that for every  $x \in \mathbb{R}$ ,  $\sin(x) = \lim_{n \rightarrow \infty} \sum_{k=1}^n (-1)^{k-1} \frac{x^{2k-1}}{(2k-1)!}$ . **Hint:** Use the formula for the error to conclude that

$$\left| \sin x - \sum_{k=1}^n (-1)^{k-1} \frac{x^{2k-1}}{(2k-1)!} \right| \leq \frac{|x|^{2n}}{(2n)!}$$

and then use the result of Problem 6.

8. Show  $\cos(x) = \lim_{n \rightarrow \infty} \sum_{k=1}^n (-1)^{k-1} \frac{x^{2k-2}}{(2k-2)!}$  by finding a suitable formula for the remainder and then using an argument similar to that done in Problem 7.
9. Using Problem 6, show  $e^x = \lim_{n \rightarrow \infty} \sum_{k=1}^n \frac{x^k}{k!}$  for all  $x \in \mathbb{R}$ .
10. Find  $a_n$  such that  $\arctan(x) = \lim_{n \rightarrow \infty} \sum_{k=0}^n a_k x^k$  for some values of  $x$ . Find the values of  $x$  for which the limit is true and prove your result. **Hint:** It is a good idea to use

$$\arctan(x) = \int_0^x \frac{1}{1+t^2} dt,$$

show

$$\frac{1}{1+t^2} = \sum_{k=0}^n (-1)^k t^{2k} \pm \frac{t^{2n+2}}{1+t^2},$$

and then integrate this finite sum from 0 to  $x$ . Thus the error would be no larger than

$$\left| \int_0^x \frac{t^{2n+2}}{1+t^2} dt \right| \leq \left| \int_0^x t^{2n+2} dt \right|.$$

11. If you did Problem 10 correctly, you found

$$\arctan x = \lim_{n \rightarrow \infty} \sum_{k=1}^n (-1)^{k-1} \frac{x^{2k-1}}{2k-1}$$

and that this limit will hold for  $x \in [-1, 1]$ . Use this to verify that

$$\frac{\pi}{4} = \lim_{n \rightarrow \infty} \sum_{k=1}^n (-1)^{k-1} \frac{1}{2k-1}.$$

12. Do for  $\ln(1+x)$  what was done for  $\arctan(x)$  and find a formula of this sort for  $\ln 2$ . Use

$$\ln(1+x) = \int_0^x \frac{1}{1+t} dt.$$

13. Repeat 10 for the function  $\ln(1-x)$ .

## 11.3 Infinite Series Of Numbers

### 11.3.1 Basic Considerations

Earlier in Definition 5.11.1 on Page 103 the notion of limit of a sequence was discussed. There is a very closely related concept called an infinite series which is dealt with in this section.

**Definition 11.3.1** *Define*

$$\sum_{k=m}^{\infty} a_k \equiv \lim_{n \rightarrow \infty} \sum_{k=m}^n a_k$$

*whenever the limit exists and is finite. In this case the series is said to converge. If it does not converge, it is said to diverge. The sequence  $\{\sum_{k=m}^n a_k\}_{n=m}^{\infty}$  in the above is called the sequence of partial sums.*

From this definition, it should be clear that infinite sums do not always make sense. Sometimes they do and sometimes they don't, depending on the behavior of the partial sums. As an example, consider  $\sum_{k=1}^{\infty} (-1)^k$ . The partial sums corresponding to this symbol alternate between  $-1$  and  $0$ . Therefore, there is no limit for the sequence of partial sums. It follows the symbol just written is meaningless and the infinite sum diverges.

**Proposition 11.3.2** *Let  $a_k \geq 0$ . Then  $\{\sum_{k=m}^n a_k\}_{n=m}^{\infty}$  is an increasing sequence. If this sequence is bounded above, then  $\sum_{k=m}^{\infty} a_k$  converges and its value equals*

$$\sup \left\{ \sum_{k=m}^n a_k : n = m, m+1, \dots \right\}.$$

*When the sequence is not bounded above,  $\sum_{k=m}^{\infty} a_k$  diverges.*

**Proof:** It follows  $\{\sum_{k=m}^n a_k\}_{n=m}^{\infty}$  is an increasing sequence because

$$\sum_{k=m}^{n+1} a_k - \sum_{k=m}^n a_k = a_{n+1} \geq 0.$$

If it is bounded above, then by the form of completeness found in Theorem 5.11.15 on Page 107 it follows the sequence of partial sums converges to  $\sup \{\sum_{k=m}^n a_k : n = m, m+1, \dots\}$ . If the sequence of partial sums is not bounded, then it is not a Cauchy sequence and so it does not converge. See Theorem 5.11.13 on Page 107. This proves the proposition.

In the case where  $a_k \geq 0$ , the above proposition shows there are only two alternatives available. Either the sequence of partial sums is bounded above or it is not bounded above. In the first case convergence occurs and in the second case, the infinite series diverges. For this reason, people will sometimes write  $\sum_{k=m}^{\infty} a_k < \infty$  to denote the case where convergence occurs and  $\sum_{k=m}^{\infty} a_k = \infty$  for the case where divergence occurs. Be very careful you never think this way in the case where it is not true that all  $a_k \geq 0$ . For example, the partial sums of  $\sum_{k=1}^{\infty} (-1)^k$  are bounded because they are all either  $-1$  or  $0$  but the series does not converge.

One of the most important examples of a convergent series is the geometric series. This series is  $\sum_{n=0}^{\infty} r^n$ . The study of this series depends on simple High school algebra and Theorem 5.11.9 on Page 106. Let  $S_n \equiv \sum_{k=0}^n r^k$ . Then

$$S_n = \sum_{k=0}^n r^k, \quad rS_n = \sum_{k=0}^n r^{k+1} = \sum_{k=1}^{n+1} r^k.$$

Therefore, subtracting the second equation from the first yields

$$(1 - r) S_n = 1 - r^{n+1}$$

and so a formula for  $S_n$  is available. In fact, if  $r \neq 1$ ,

$$S_n = \frac{1 - r^{n+1}}{1 - r}.$$

By Theorem 5.11.9,  $\lim_{n \rightarrow \infty} S_n = \frac{1}{1-r}$  in the case when  $|r| < 1$ . Now if  $|r| \geq 1$ , the limit clearly does not exist because  $S_n$  fails to be a Cauchy sequence (Why?). This shows the following.

**Theorem 11.3.3** *The geometric series,  $\sum_{n=0}^{\infty} r^n$  converges and equals  $\frac{1}{1-r}$  if  $|r| < 1$  and diverges if  $|r| \geq 1$ .*

If the series do converge, the following holds.

**Theorem 11.3.4** *If  $\sum_{k=m}^{\infty} a_k$  and  $\sum_{k=m}^{\infty} b_k$  both converge and  $x, y$  are numbers, then*

$$\sum_{k=m}^{\infty} a_k = \sum_{k=m+j}^{\infty} a_{k-j} \quad (11.3)$$

$$\sum_{k=m}^{\infty} xa_k + yb_k = x \sum_{k=m}^{\infty} a_k + y \sum_{k=m}^{\infty} b_k \quad (11.4)$$

$$\left| \sum_{k=m}^{\infty} a_k \right| \leq \sum_{k=m}^{\infty} |a_k| \quad (11.5)$$

where in the last inequality, the last sum equals  $+\infty$  if the partial sums are not bounded above.

**Proof:** The above theorem is really only a restatement of Theorem 5.11.6 on Page 104 and the above definitions of infinite series. Thus

$$\sum_{k=m}^{\infty} a_k = \lim_{n \rightarrow \infty} \sum_{k=m}^n a_k = \lim_{n \rightarrow \infty} \sum_{k=m+j}^{n+j} a_{k-j} = \sum_{k=m+j}^{\infty} a_{k-j}.$$

To establish (11.4), use Theorem 5.11.6 on Page 104 to write

$$\begin{aligned} \sum_{k=m}^{\infty} xa_k + yb_k &= \lim_{n \rightarrow \infty} \sum_{k=m}^n xa_k + yb_k \\ &= \lim_{n \rightarrow \infty} \left( x \sum_{k=m}^n a_k + y \sum_{k=m}^n b_k \right) \\ &= x \sum_{k=m}^{\infty} a_k + y \sum_{k=m}^{\infty} b_k. \end{aligned}$$

Formula (11.5) follows from the observation that, from the triangle inequality,

$$\left| \sum_{k=m}^n a_k \right| \leq \sum_{k=m}^n |a_k|$$

and so

$$\left| \sum_{k=m}^{\infty} a_k \right| = \lim_{n \rightarrow \infty} \left| \sum_{k=m}^n a_k \right| \leq \sum_{k=m}^{\infty} |a_k|.$$

**Example 11.3.5** Find  $\sum_{n=0}^{\infty} \left( \frac{5}{2^n} + \frac{6}{3^n} \right)$ .

From the above theorem and Theorem 11.3.3,

$$\begin{aligned} \sum_{n=0}^{\infty} \left( \frac{5}{2^n} + \frac{6}{3^n} \right) &= 5 \sum_{n=0}^{\infty} \frac{1}{2^n} + 6 \sum_{n=0}^{\infty} \frac{1}{3^n} \\ &= 5 \frac{1}{1 - (1/2)} + 6 \frac{1}{1 - (1/3)} = 19. \end{aligned}$$

The following criterion is useful in checking convergence.

**Theorem 11.3.6** The sum  $\sum_{k=m}^{\infty} a_k$  converges if and only if for all  $\varepsilon > 0$ , there exists  $n_\varepsilon$  such that if  $q \geq p \geq n_\varepsilon$ , then

$$\left| \sum_{k=p}^q a_k \right| < \varepsilon. \quad (11.6)$$

**Proof:** Suppose first that the series converges. Then  $\{\sum_{k=m}^n a_k\}_{n=m}^{\infty}$  is a Cauchy sequence by Theorem 5.11.13 on Page 107. Therefore, there exists  $n_\varepsilon > m$  such that if  $q \geq p - 1 \geq n_\varepsilon > m$ ,

$$\left| \sum_{k=m}^q a_k - \sum_{k=m}^{p-1} a_k \right| = \left| \sum_{k=p}^q a_k \right| < \varepsilon. \quad (11.7)$$

Next suppose (11.6) holds. Then from (11.7) it follows upon letting  $p$  be replaced with  $p + 1$  that  $\{\sum_{k=m}^n a_k\}_{n=m}^{\infty}$  is a Cauchy sequence and so, by the completeness axiom, it converges. By the definition of infinite series, this shows the infinite sum converges as claimed.

**Definition 11.3.7** A series

$$\sum_{k=m}^{\infty} a_k$$

is said to converge absolutely if

$$\sum_{k=m}^{\infty} |a_k|$$

converges. If the series does converge but does not converge absolutely, then it is said to converge conditionally.

**Theorem 11.3.8** If  $\sum_{k=m}^{\infty} a_k$  converges absolutely, then it converges.

**Proof:** Let  $\varepsilon > 0$  be given. Then by assumption and Theorem 11.3.6, there exists  $n_\varepsilon$  such that whenever  $q \geq p \geq n_\varepsilon$ ,

$$\sum_{k=p}^q |a_k| < \varepsilon.$$

Therefore, from the triangle inequality,

$$\varepsilon > \sum_{k=p}^q |a_k| \geq \left| \sum_{k=p}^q a_k \right|.$$

By Theorem 11.3.6,  $\sum_{k=m}^{\infty} a_k$  converges and this proves the theorem.

In fact, the above theorem is really another version of the completeness axiom. Thus its validity implies completeness. You might try to show this.

**Theorem 11.3.9** (comparison test) Suppose  $\{a_n\}$  and  $\{b_n\}$  are sequences of non negative real numbers and suppose for all  $n$  sufficiently large,  $a_n \leq b_n$ . Then

1. If  $\sum_{n=k}^{\infty} b_n$  converges, then  $\sum_{n=m}^{\infty} a_n$  converges.
2. If  $\sum_{n=k}^{\infty} a_n$  diverges, then  $\sum_{n=m}^{\infty} b_n$  diverges.

**Proof:** Consider the first claim. From the assumption there exists  $n^*$  such that  $n^* > \max(k, m)$  and for all  $n \geq n^*$   $b_n \geq a_n$ . Then if  $p \geq n^*$ ,

$$\begin{aligned} \sum_{n=m}^p a_n &\leq \sum_{n=m}^{n^*} a_n + \sum_{n=n^*+1}^k b_n \\ &\leq \sum_{n=m}^{n^*} a_n + \sum_{n=k}^{\infty} b_n. \end{aligned}$$

Thus the sequence,  $\{\sum_{n=m}^p a_n\}_{p=m}^{\infty}$  is bounded above and increasing. Therefore, it converges by completeness. The second claim is left as an exercise.

**Example 11.3.10** Determine the convergence of  $\sum_{n=1}^{\infty} \frac{1}{n^2}$ .

For  $n > 1$ ,

$$\frac{1}{n^2} \leq \frac{1}{n(n-1)}.$$

Now

$$\begin{aligned} \sum_{n=2}^p \frac{1}{n(n-1)} &= \sum_{n=2}^p \left[ \frac{1}{n-1} - \frac{1}{n} \right] \\ &= 1 - \frac{1}{p} \rightarrow 1 \end{aligned}$$

Therefore, letting  $a_n = \frac{1}{n^2}$  and  $b_n = \frac{1}{n(n-1)}$

A convenient way to implement the comparison test is to use the limit comparison test. This is considered next.

**Theorem 11.3.11** Let  $a_n, b_n > 0$  and suppose for all  $n$  large enough,

$$0 < a < \frac{a_n}{b_n} \leq \frac{a_n}{b_n} < b < \infty.$$

Then  $\sum a_n$  and  $\sum b_n$  converge or diverge together.

**Proof:** Let  $n^*$  be such that  $n \geq n^*$ , then

$$\frac{a_n}{b_n} > a \text{ and } \frac{a_n}{b_n} < b$$

and so for all such  $n$ ,

$$ab_n < a_n < bb_n$$

and so the conclusion follows from the comparison test.

**Example 11.3.12** Determine the convergence of  $\sum_{k=1}^{\infty} \frac{1}{\sqrt{n^4+2n+7}}$ .



This series converges by the limit comparison test. Compare with the series of Example 11.3.10.

$$\begin{aligned}\lim_{n \rightarrow \infty} \frac{\left(\frac{1}{n^2}\right)}{\left(\frac{1}{\sqrt{n^4+2n+7}}\right)} &= \lim_{n \rightarrow \infty} \frac{\sqrt{n^4+2n+7}}{n^2} \\ &= \lim_{n \rightarrow \infty} \sqrt{1 + \frac{2}{n^3} + \frac{7}{n^4}} = 1.\end{aligned}$$

Therefore, the series converges with the series of Example 11.3.10. How did I know what to compare with? I noticed that  $\sqrt{n^4+2n+7}$  is essentially like  $\sqrt{n^4} = n^2$  for large enough  $n$ . You see, the higher order term,  $n^4$  dominates the other terms in  $n^4+2n+7$ . Therefore, reasoning that  $1/\sqrt{n^4+2n+7}$  is a lot like  $1/n^2$  for large  $n$ , it was easy to see what to compare with. Of course this is not always easy and there is room for acquiring skill through practice.

To really exploit this limit comparison test, it is desirable to get lots of examples of series, some which converge and some which do not. The tool for obtaining these examples here will be the following wonderful theorem known as the Cauchy condensation test.

**Theorem 11.3.13** *Let  $a_n \geq 0$  and suppose the terms of the sequence  $\{a_n\}$  are decreasing. Thus  $a_n \geq a_{n+1}$  for all  $n$ . Then*

$$\sum_{n=1}^{\infty} a_n \text{ and } \sum_{n=0}^{\infty} 2^n a_{2^n}$$

*converge or diverge together.*

**Proof:** This follows from the inequality of the following claim.

**Claim:**

$$\sum_{k=1}^n 2^k a_{2^{k-1}} \geq \sum_{k=1}^{2^n} a_k \geq \sum_{k=0}^n 2^{k-1} a_{2^k}.$$

**Proof of the Claim:** Note the claim is true for  $n = 1$ . Suppose the claim is true for  $n$ . Then, since  $2^{n+1} - 2^n = 2^n$ , and the terms,  $a_n$ , are decreasing,

$$\begin{aligned}\sum_{k=1}^{n+1} 2^k a_{2^{k-1}} &= 2^{n+1} a_{2^n} + \sum_{k=1}^n 2^k a_{2^{k-1}} \geq 2^{n+1} a_{2^n} + \sum_{k=1}^{2^n} a_k \\ &\geq \sum_{k=1}^{2^{n+1}} a_k \geq 2^n a_{2^{n+1}} + \sum_{k=1}^{2^n} a_k \geq 2^n a_{2^{n+1}} + \sum_{k=0}^n 2^{k-1} a_{2^k} = \sum_{k=0}^{n+1} 2^{k-1} a_{2^k}.\end{aligned}$$

**Example 11.3.14** *Determine the convergence of  $\sum_{k=1}^{\infty} \frac{1}{k^p}$  where  $p$  is a positive number. These are called the  $p$  series.*

Let  $a_n = \frac{1}{n^p}$ . Then  $a_{2^n} = \left(\frac{1}{2^p}\right)^n$ . From the Cauchy condensation test the two series

$$\sum_{n=1}^{\infty} \frac{1}{n^p} \text{ and } \sum_{n=0}^{\infty} 2^n \left(\frac{1}{2^p}\right)^n = \sum_{n=0}^{\infty} \left(2^{(1-p)}\right)^n$$

converge or diverge together. If  $p > 1$ , the last series above is a geometric series having common ratio less than 1 and so it converges. If  $p \leq 1$ , it is still a geometric series but in this case the common ratio is either 1 or greater than 1 so the series diverges. It follows that the  $p$  series converges if  $p > 1$  and diverges if  $p \leq 1$ . In particular,  $\sum_{n=1}^{\infty} n^{-1}$  diverges while  $\sum_{n=1}^{\infty} n^{-2}$  converges.

**Example 11.3.15** Determine the convergence of  $\sum_{k=1}^{\infty} \frac{1}{\sqrt{n^2+100n}}$ .

Use the limit comparison test.

$$\lim_{n \rightarrow \infty} \frac{\left(\frac{1}{n}\right)}{\left(\frac{1}{\sqrt{n^2+100n}}\right)} = 1$$

and so this series diverges with  $\sum_{k=1}^{\infty} \frac{1}{k}$ .

**Example 11.3.16** Determine the convergence of  $\sum_{k=2}^{\infty} \frac{1}{k \ln k}$ .

Use the Cauchy condensation test. The above series does the same thing in terms of convergence as the series

$$\sum_{n=1}^{\infty} 2^n \frac{1}{2^n \ln(2^n)} = \sum_{n=1}^{\infty} \frac{1}{n \ln 2}$$

and this series diverges by limit comparison with the series  $\sum \frac{1}{n}$ .

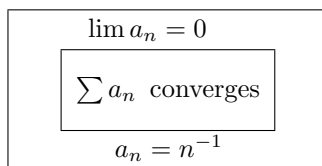
Sometimes it is good to be able to say a series does not converge. The  $n^{\text{th}}$  term test gives such a condition which is sufficient for this. It is really a corollary of Theorem 11.3.6.

**Theorem 11.3.17** If  $\sum_{n=m}^{\infty} a_n$  converges, then  $\lim_{n \rightarrow \infty} a_n = 0$ .

**Proof:** Apply Theorem 11.3.6 to conclude that

$$\lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} \sum_{k=n}^n a_k = 0.$$

It is very important to observe that this theorem goes only in one direction. That is, you cannot conclude the series converges if  $\lim_{n \rightarrow \infty} a_n = 0$ . If this happens, you don't know anything from this information. Recall  $\lim_{n \rightarrow \infty} n^{-1} = 0$  but  $\sum_{n=1}^{\infty} n^{-1}$  diverges. The following picture is descriptive of the situation.



### 11.3.2 More Tests For Convergence

So far, the tests for convergence have been applied to non negative terms only. Sometimes, a series converges, not because the terms of the series get small fast enough, but because of cancellation taking place between positive and negative terms. A discussion of this involves some simple algebra.

Let  $\{a_n\}$  and  $\{b_n\}$  be sequences and let

$$A_n \equiv \sum_{k=1}^n a_k, \quad A_{-1} \equiv A_0 \equiv 0.$$

Then if  $p < q$

$$\begin{aligned} \sum_{n=p}^q a_n b_n &= \sum_{n=p}^q b_n (A_n - A_{n-1}) = \sum_{n=p}^q b_n A_n - \sum_{n=p}^q b_n A_{n-1} \\ &= \sum_{n=p}^q b_n A_n - \sum_{n=p-1}^{q-1} b_{n+1} A_n = b_q A_q - b_p A_{p-1} + \sum_{n=p}^{q-1} A_n (b_n - b_{n+1}) \end{aligned}$$

This formula is called the partial summation formula. It is just like integration by parts.

**Theorem 11.3.18** (*Dirichlet's test*) Suppose  $A_n$  is bounded and  $\lim_{n \rightarrow \infty} b_n = 0$ , with  $b_n \geq b_{n+1}$ . Then

$$\sum_{n=1}^{\infty} a_n b_n$$

converges.

**Proof:** This follows quickly from Theorem 11.3.6. Indeed, letting  $|A_n| \leq C$ , and using the partial summation formula above along with the assumption that the  $b_n$  are decreasing,

$$\begin{aligned} \left| \sum_{n=p}^q a_n b_n \right| &= \left| b_q A_q - b_p A_{p-1} + \sum_{n=p}^{q-1} A_n (b_n - b_{n+1}) \right| \\ &\leq C(|b_q| + |b_p|) + C \sum_{n=p}^{q-1} (b_n - b_{n+1}) \\ &= C(|b_q| + |b_p|) + C(b_p - b_q) \end{aligned}$$

and by assumption, this last expression is small whenever  $p$  and  $q$  are sufficiently large. This proves the theorem.

The following corollary is known as the alternating series test.

**Corollary 11.3.19** If  $\lim_{n \rightarrow \infty} b_n = 0$ , with  $b_n \geq b_{n+1}$ , then  $\sum_{n=1}^{\infty} (-1)^n b_n$  converges.

A favorite test for convergence is the ratio test. This is discussed next.

**Theorem 11.3.20** Suppose  $|a_n| > 0$  for all  $n$  and suppose

$$\lim_{n \rightarrow \infty} \frac{|a_{n+1}|}{|a_n|} = r.$$

Then

$$\sum_{n=1}^{\infty} a_n \begin{cases} \text{diverges if } r > 1 \\ \text{converges absolutely if } r < 1 \\ \text{test fails if } r = 1 \end{cases}.$$

**Proof:** Suppose  $r < 1$ . Then there exists  $n_1$  such that if  $n \geq n_1$ , then

$$0 < \left| \frac{a_{n+1}}{a_n} \right| < R$$

where  $r < R < 1$ . Then

$$|a_{n+1}| < R |a_n|$$

for all such  $n$ . Therefore,

$$|a_{n_1+p}| < R |a_{n_1+p-1}| < R^2 |a_{n_1+p-2}| < \cdots < R^p |a_{n_1}| \quad (11.8)$$

and so if  $m > n$ , then  $|a_m| < R^{m-n_1} |a_{n_1}|$ . By the comparison test and the theorem on geometric series,  $\sum |a_n|$  converges. This proves the convergence part of the theorem.

To verify the divergence part, note that if  $r > 1$ , then (11.8) can be turned around for some  $R > 1$ . Showing  $\lim_{n \rightarrow \infty} |a_n| = \infty$ . Since the  $n^{\text{th}}$  term fails to converge to 0, it follows the series diverges.

To see the test fails if  $r = 1$ , consider  $\sum n^{-1}$  and  $\sum n^{-2}$ . The first series diverges while the second one converges but in both cases,  $r = 1$ . (Be sure to check this last claim.)

The ratio test is very useful for many different examples but it is somewhat unsatisfactory mathematically. One reason for this is the assumption that  $a_n > 0$ , necessitated by the need to divide by  $a_n$ , and the other reason is the possibility that the limit might not exist. The next test, called the root test removes both of these objections.

**Theorem 11.3.21** Suppose  $|a_n|^{1/n} < R < 1$  for all  $n$  sufficiently large. Then

$$\sum_{n=1}^{\infty} a_n \text{ converges absolutely.}$$

If there are infinitely many values of  $n$  such that  $|a_n|^{1/n} \geq 1$ , then

$$\sum_{n=1}^{\infty} a_n \text{ diverges.}$$

**Proof:** Suppose first that  $|a_n|^{1/n} < R < 1$  for all  $n$  sufficiently large. Say this holds for all  $n \geq n_R$ . Then for such  $n$ ,

$$\sqrt[n]{|a_n|} < R.$$

Therefore, for such  $n$ ,

$$|a_n| \leq R^n$$

and so the comparison test with a geometric series applies and gives absolute convergence as claimed.

Next suppose  $|a_n|^{1/n} \geq 1$  for infinitely many values of  $n$ . Then for those values of  $n$ ,  $|a_n| \geq 1$  and so the series fails to converge by the  $n^{\text{th}}$  term test.

**Corollary 11.3.22** Suppose  $\lim_{n \rightarrow \infty} |a_n|^{1/n}$  exists and equals  $r$ . Then

$$\sum_{k=m}^{\infty} a_k \begin{cases} \text{converges absolutely if } r < 1 \\ \text{test fails if } r = 1 \\ \text{diverges if } r > 1 \end{cases}$$

**Proof:** The first and last alternatives follow from Theorem 11.3.21. To see the test fails if  $r = 1$ , consider the two series  $\sum_{n=1}^{\infty} \frac{1}{n}$  and  $\sum_{n=1}^{\infty} \frac{1}{n^2}$  both of which have  $r = 1$  but having different convergence properties.

**Example 11.3.23** Show that for all  $x \in \mathbb{R}$ ,

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}. \quad (11.9)$$

By Taylor's theorem

$$\begin{aligned} e^x &= 1 + x + \frac{x^2}{2!} + \cdots + \frac{x^n}{n!} + e^{\xi_n} \frac{x^{n+1}}{(n+1)!} \\ &= \sum_{k=0}^n \frac{x^k}{k!} + e^{\xi_n} \frac{x^{n+1}}{(n+1)!} \end{aligned} \quad (11.10)$$

where  $|\xi_n| \leq |x|$ . Now for any  $x \in \mathbb{R}$

$$\left| e^{\xi_n} \frac{x^{n+1}}{(n+1)!} \right| \leq e^{|x|} \frac{|x|^{n+1}}{(n+1)!}$$

and an application of the ratio test shows

$$\sum_{n=0}^{\infty} e^{|x|} \frac{|x|^{n+1}}{(n+1)!} < \infty.$$

Therefore, the  $n^{\text{th}}$  term converges to zero by the  $n^{\text{th}}$  term test and so for each  $x \in \mathbb{R}$ ,

$$\lim_{n \rightarrow \infty} e^{\xi_n} \frac{x^{n+1}}{(n+1)!} = 0.$$

Therefore, taking the limit in (11.10) it follows (11.9) holds.

### 11.3.3 Double Series

Sometimes it is required to consider double series which are of the form

$$\sum_{k=m}^{\infty} \sum_{j=m}^{\infty} a_{jk} \equiv \sum_{k=m}^{\infty} \left( \sum_{j=m}^{\infty} a_{jk} \right).$$

In other words, first sum on  $j$  yielding something which depends on  $k$  and then sum these. The major consideration for these double series is the question of when

$$\sum_{k=m}^{\infty} \sum_{j=m}^{\infty} a_{jk} = \sum_{j=m}^{\infty} \sum_{k=m}^{\infty} a_{jk}.$$

In other words, when does it make no difference which subscript is summed over first? In the case of finite sums there is no issue here. You can always write

$$\sum_{k=m}^M \sum_{j=m}^N a_{jk} = \sum_{j=m}^N \sum_{k=m}^M a_{jk}$$

because addition is commutative. However, there are limits involved with infinite sums and the interchange in order of summation involves taking limits in a different order. Therefore, it is not always true that it is permissible to interchange the two sums. A general rule of thumb is this: If something involves changing the order in which two limits are taken, you may not do it without agonizing over the question. In general, limits foul up algebra and also introduce things which are counter intuitive. Here is an example. This example is a little technical. It is placed here just to prove conclusively there is a question which needs to be considered.

**Example 11.3.24** Consider the following picture which depicts some of the ordered pairs  $(m, n)$  where  $m, n$  are positive integers.

$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$c_{\bullet}$	$0_{\bullet}$	$-c_{\bullet}$
$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$c_{\bullet}$	$0_{\bullet}$	$-c_{\bullet}$	$0_{\bullet}$
$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$c_{\bullet}$	$0_{\bullet}$	$-c_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$
$0_{\bullet}$	$0_{\bullet}$	$c_{\bullet}$	$0_{\bullet}$	$-c_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$
$0_{\bullet}$	$c_{\bullet}$	$0_{\bullet}$	$-c_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$
$b_{\bullet}$	$0_{\bullet}$	$-c_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$
$0_{\bullet}$	$a_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$	$0_{\bullet}$

The numbers next to the point are the values of  $a_{mn}$ . You see  $a_{nn} = 0$  for all  $n$ ,  $a_{21} = a$ ,  $a_{12} = b$ ,  $a_{mn} = c$  for  $(m, n)$  on the line  $y = 1 + x$  whenever  $m > 1$ , and  $a_{mn} = -c$  for all  $(m, n)$  on the line  $y = x - 1$  whenever  $m > 2$ .

Then  $\sum_{m=1}^{\infty} a_{mn} = a$  if  $n = 1$ ,  $\sum_{m=1}^{\infty} a_{mn} = b - c$  if  $n = 2$  and if  $n > 2$ ,  $\sum_{m=1}^{\infty} a_{mn} = 0$ . Therefore,

$$\sum_{n=1}^{\infty} \sum_{m=1}^{\infty} a_{mn} = a + b - c.$$

Next observe that  $\sum_{n=1}^{\infty} a_{mn} = b$  if  $m = 1$ ,  $\sum_{n=1}^{\infty} a_{mn} = a + c$  if  $m = 2$ , and  $\sum_{n=1}^{\infty} a_{mn} = 0$  if  $m > 2$ . Therefore,

$$\sum_{m=1}^{\infty} \sum_{n=1}^{\infty} a_{mn} = b + a + c$$

and so the two sums are different. Moreover, you can see that by assigning different values of  $a, b$ , and  $c$ , you can get an example for any two different numbers desired.

Don't become upset by this. It happens because, as indicated above, limits are taken in two different orders. An infinite sum always involves a limit and this illustrates why you must always remember this. This example in no way violates the commutative law of addition which has nothing to do with limits. However, it turns out that if  $a_{ij} \geq 0$  for all  $i, j$ , then you can always interchange the order of summation. This is shown next and is based on the following lemma. First, some notation should be discussed.

**Definition 11.3.25** Let  $f(a, b) \in [-\infty, \infty]$  for  $a \in A$  and  $b \in B$  where  $A, B$  are sets which means that  $f(a, b)$  is either a number,  $\infty$ , or  $-\infty$ . The symbol,  $+\infty$  is interpreted as a point out at the end of the number line which is larger than every real number. Of course there is no such number. That is why it is called  $\infty$ . The symbol,  $-\infty$  is interpreted similarly. Then  $\sup_{a \in A} f(a, b)$  means  $\sup(S_b)$  where  $S_b \equiv \{f(a, b) : a \in A\}$ .

Unlike limits, you can take the sup in different orders.

**Lemma 11.3.26** Let  $f(a, b) \in [-\infty, \infty]$  for  $a \in A$  and  $b \in B$  where  $A, B$  are sets. Then

$$\sup_{a \in A} \sup_{b \in B} f(a, b) = \sup_{b \in B} \sup_{a \in A} f(a, b).$$

**Proof:** Note that for all  $a, b$ ,  $f(a, b) \leq \sup_{b \in B} \sup_{a \in A} f(a, b)$  and therefore, for all  $a$ ,  $\sup_{b \in B} f(a, b) \leq \sup_{b \in B} \sup_{a \in A} f(a, b)$ . Therefore,

$$\sup_{a \in A} \sup_{b \in B} f(a, b) \leq \sup_{b \in B} \sup_{a \in A} f(a, b).$$

Repeat the same argument interchanging  $a$  and  $b$ , to get the conclusion of the lemma.

**Lemma 11.3.27** If  $\{A_n\}$  is an increasing sequence in  $[-\infty, \infty]$ , then  $\sup\{A_n\} = \lim_{n \rightarrow \infty} A_n$ .

**Proof:** Let  $\sup\{A_n : n \in \mathbb{N}\} = r$ . In the first case, suppose  $r < \infty$ . Then letting  $\varepsilon > 0$  be given, there exists  $n$  such that  $A_n \in (r - \varepsilon, r]$ . Since  $\{A_n\}$  is increasing, it follows if  $m > n$ , then  $r - \varepsilon < A_n \leq A_m \leq r$  and so  $\lim_{n \rightarrow \infty} A_n = r$  as claimed. In the case where  $r = \infty$ , then if  $a$  is a real number, there exists  $n$  such that  $A_n > a$ . Since  $\{A_k\}$  is increasing, it follows that if  $m > n$ ,  $A_m > a$ . But this is what is meant by  $\lim_{n \rightarrow \infty} A_n = \infty$ . The other case is that  $r = -\infty$ . But in this case,  $A_n = -\infty$  for all  $n$  and so  $\lim_{n \rightarrow \infty} A_n = -\infty$ .

**Theorem 11.3.28** Let  $a_{ij} \geq 0$ . Then  $\sum_{i=1}^{\infty} \sum_{j=1}^{\infty} a_{ij} = \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} a_{ij}$ .

**Proof:** First note there is no trouble in defining these sums because the  $a_{ij}$  are all nonnegative. If a sum diverges, it only diverges to  $\infty$  and so  $\infty$  is the value of the sum. Next note that

$$\sum_{j=r}^{\infty} \sum_{i=r}^{\infty} a_{ij} \geq \sup_n \sum_{j=r}^{\infty} \sum_{i=r}^n a_{ij}$$

because for all  $j$ ,

$$\sum_{i=r}^{\infty} a_{ij} \geq \sum_{i=r}^n a_{ij}.$$

Therefore,

$$\begin{aligned} \sum_{j=r}^{\infty} \sum_{i=r}^{\infty} a_{ij} &\geq \sup_n \sum_{j=r}^{\infty} \sum_{i=r}^n a_{ij} = \sup_n \lim_{m \rightarrow \infty} \sum_{j=r}^m \sum_{i=r}^n a_{ij} \\ &= \sup_n \lim_{m \rightarrow \infty} \sum_{i=r}^n \sum_{j=r}^m a_{ij} = \sup_n \sum_{i=r}^n \lim_{m \rightarrow \infty} \sum_{j=r}^m a_{ij} \\ &= \sup_n \sum_{i=r}^n \sum_{j=r}^{\infty} a_{ij} = \lim_{n \rightarrow \infty} \sum_{i=r}^n \sum_{j=r}^{\infty} a_{ij} = \sum_{i=r}^{\infty} \sum_{j=r}^{\infty} a_{ij} \end{aligned}$$

Interchanging the  $i$  and  $j$  in the above argument proves the theorem.

The following is the fundamental result on double sums.

**Theorem 11.3.29** Let  $a_{ij}$  be a number and suppose

$$\sum_{i=r}^{\infty} \sum_{j=r}^{\infty} |a_{ij}| < \infty.$$

Then

$$\sum_{i=r}^{\infty} \sum_{j=r}^{\infty} a_{ij} = \sum_{j=r}^{\infty} \sum_{i=r}^{\infty} a_{ij}$$

and every infinite sum encountered in the above equation converges.

**Proof:** By Theorem 11.3.28

$$\sum_{j=r}^{\infty} \sum_{i=r}^{\infty} |a_{ij}| = \sum_{i=r}^{\infty} \sum_{j=r}^{\infty} |a_{ij}| < \infty$$

Therefore, for each  $j$ ,  $\sum_{i=r}^{\infty} |a_{ij}| < \infty$  and for each  $i$ ,  $\sum_{j=r}^{\infty} |a_{ij}| < \infty$ . By Theorem 11.3.8 on Page 263,  $\sum_{i=r}^{\infty} a_{ij}$ ,  $\sum_{j=r}^{\infty} a_{ij}$  both converge, the first one for every  $j$  and the second for every  $i$ . Also,

$$\sum_{j=r}^{\infty} \left| \sum_{i=r}^{\infty} a_{ij} \right| \leq \sum_{j=r}^{\infty} \sum_{i=r}^{\infty} |a_{ij}| < \infty$$

and

$$\sum_{i=r}^{\infty} \left| \sum_{j=r}^{\infty} a_{ij} \right| \leq \sum_{i=r}^{\infty} \sum_{j=r}^{\infty} |a_{ij}| < \infty$$

so by Theorem 11.3.8 again,

$$\sum_{j=r}^{\infty} \sum_{i=r}^{\infty} a_{ij}, \sum_{i=r}^{\infty} \sum_{j=r}^{\infty} a_{ij}$$

both exist. It only remains to verify they are equal. Note  $0 \leq (|a_{ij}| + a_{ij}) \leq |a_{ij}|$ . Therefore, by Theorem 11.3.28 and Theorem 11.3.4 on Page 262

$$\begin{aligned} \sum_{j=r}^{\infty} \sum_{i=r}^{\infty} |a_{ij}| + \sum_{j=r}^{\infty} \sum_{i=r}^{\infty} a_{ij} &= \sum_{j=r}^{\infty} \sum_{i=r}^{\infty} (|a_{ij}| + a_{ij}) \\ &= \sum_{i=r}^{\infty} \sum_{j=r}^{\infty} (|a_{ij}| + a_{ij}) \\ &= \sum_{i=r}^{\infty} \sum_{j=r}^{\infty} |a_{ij}| + \sum_{i=r}^{\infty} \sum_{j=r}^{\infty} a_{ij} \\ &= \sum_{j=r}^{\infty} \sum_{i=r}^{\infty} |a_{ij}| + \sum_{i=r}^{\infty} \sum_{j=r}^{\infty} a_{ij} \end{aligned}$$

and so  $\sum_{j=r}^{\infty} \sum_{i=r}^{\infty} a_{ij} = \sum_{i=r}^{\infty} \sum_{j=r}^{\infty} a_{ij}$  as claimed. This proves the theorem.

One of the most important applications of this theorem is to the problem of multiplication of series.

**Definition 11.3.30** Let  $\sum_{i=r}^{\infty} a_i$  and  $\sum_{i=r}^{\infty} b_i$  be two series. For  $n \geq r$ , define

$$c_n \equiv \sum_{k=r}^n a_k b_{n-k+r}.$$

The series  $\sum_{n=r}^{\infty} c_n$  is called the Cauchy product of the two series.



It isn't hard to see where this comes from. Formally write the following in the case  $r = 0$ :

$$(a_0 + a_1 + a_2 + a_3 \cdots)(b_0 + b_1 + b_2 + b_3 \cdots)$$

and start multiplying in the usual way. This yields

$$a_0b_0 + (a_0b_1 + b_0a_1) + (a_0b_2 + a_1b_1 + a_2b_0) + \cdots$$

and you see the expressions in parentheses above are just the  $c_n$  for  $n = 0, 1, 2, \dots$ . Therefore, it is reasonable to conjecture that

$$\sum_{i=r}^{\infty} a_i \sum_{j=r}^{\infty} b_j = \sum_{n=r}^{\infty} c_n$$

and of course there would be no problem with this in the case of finite sums but in the case of infinite sums, it is necessary to prove a theorem. The following is a special case of Merten's theorem.

**Theorem 11.3.31** Suppose  $\sum_{i=r}^{\infty} a_i$  and  $\sum_{j=r}^{\infty} b_j$  both converge absolutely<sup>1</sup>. Then

$$\left( \sum_{i=r}^{\infty} a_i \right) \left( \sum_{j=r}^{\infty} b_j \right) = \sum_{n=r}^{\infty} c_n$$

where

$$c_n = \sum_{k=r}^n a_k b_{n-k+r}.$$

**Proof:** Let  $p_{nk} = 1$  if  $r \leq k \leq n$  and  $p_{nk} = 0$  if  $k > n$ . Then

$$c_n = \sum_{k=r}^{\infty} p_{nk} a_k b_{n-k+r}.$$

Also,

$$\begin{aligned} \sum_{k=r}^{\infty} \sum_{n=r}^{\infty} p_{nk} |a_k| |b_{n-k+r}| &= \sum_{k=r}^{\infty} |a_k| \sum_{n=r}^{\infty} p_{nk} |b_{n-k+r}| \\ &= \sum_{k=r}^{\infty} |a_k| \sum_{n=k}^{\infty} |b_{n-k+r}| \\ &= \sum_{k=r}^{\infty} |a_k| \sum_{n=k}^{\infty} |b_{n-(k-r)}| \\ &= \sum_{k=r}^{\infty} |a_k| \sum_{m=r}^{\infty} |b_m| < \infty. \end{aligned}$$

Therefore, by Theorem 11.3.29

$$\begin{aligned} \sum_{n=r}^{\infty} c_n &= \sum_{n=r}^{\infty} \sum_{k=r}^n a_k b_{n-k+r} = \sum_{n=r}^{\infty} \sum_{k=r}^{\infty} p_{nk} a_k b_{n-k+r} \\ &= \sum_{k=r}^{\infty} a_k \sum_{n=r}^{\infty} p_{nk} b_{n-k+r} = \sum_{k=r}^{\infty} a_k \sum_{n=k}^{\infty} b_{n-k+r} \\ &= \sum_{k=r}^{\infty} a_k \sum_{m=r}^{\infty} b_m \end{aligned}$$

<sup>1</sup>Actually, it is only necessary to assume one of the series converges and the other converges absolutely. This is known as Merten's theorem and may be read in the 1974 book by Apostol listed in the bibliography.

and this proves the theorem.

## 11.4 Exercises

1. Determine whether the following series converge absolutely, conditionally, or not at all and give reasons for your answers.

(a)  $\sum_{n=1}^{\infty} (-1)^n \frac{1}{\sqrt{n^2+n+1}}$

(b)  $\sum_{n=1}^{\infty} (-1)^n (\sqrt{n+1} - \sqrt{n})$

(c)  $\sum_{n=1}^{\infty} (-1)^n \frac{(n!)^2}{(2n)!}$

(d)  $\sum_{n=1}^{\infty} (-1)^n \frac{(2n)!}{(n!)^2}$

(e)  $\sum_{n=1}^{\infty} \frac{(-1)^n}{2n+2}$

(f)  $\sum_{n=1}^{\infty} (-1)^n \left(\frac{n}{n+1}\right)^n$

(g)  $\sum_{n=1}^{\infty} (-1)^n \left(\frac{n}{n+1}\right)^{n^2}$

2. Determine whether the following series converge absolutely, conditionally, or not at all and give reasons for your answers.

(a)  $\sum_{n=1}^{\infty} (-1)^n \frac{\ln(k^5)}{k}$

(b)  $\sum_{n=1}^{\infty} (-1)^n \frac{\ln(k^5)}{k^{1.01}}$

(c)  $\sum_{n=1}^{\infty} (-1)^n \frac{10^n}{(1.01)^n}$

(d)  $\sum_{n=1}^{\infty} (-1)^n \sin\left(\frac{1}{n}\right)$

(e)  $\sum_{n=1}^{\infty} (-1)^n \tan\left(\frac{1}{n^2}\right)$

(f)  $\sum_{n=1}^{\infty} (-1)^n \cos\left(\frac{1}{n^2}\right)$

3. Suppose  $\sum_{n=1}^{\infty} a_n$  converges absolutely, can the same thing be said about  $\sum_{n=1}^{\infty} a_n^2$ ? Explain.
4. A person says a series converges conditionally by the ratio test. Explain why his statement is total nonsense.
5. A person says a series diverges by the alternating series test. Explain why his statement is total nonsense.
6. Find a series which diverges using one test but converges using another if possible. If this is not possible, tell why.
7. If  $\sum_{n=1}^{\infty} a_n$  and  $\sum_{n=1}^{\infty} b_n$  both converge, can you conclude the sum,  $\sum_{n=1}^{\infty} a_n b_n$  converges?
8. If  $\sum_{n=1}^{\infty} a_n$  converges absolutely, and  $b_n$  is bounded, can you conclude  $\sum_{n=1}^{\infty} a_n b_n$  converges? What if it is only the case that  $\sum_{n=1}^{\infty} a_n$  converges?
9. The logarithm test states the following. Suppose  $a_k \neq 0$  for large  $k$  and that  $p = \lim_{k \rightarrow \infty} \frac{\ln\left(\frac{1}{|a_k|}\right)}{\ln k}$  exists. If  $p > 1$ , then  $\sum_{k=1}^{\infty} a_k$  converges absolutely. If  $p < 1$ , then the series,  $\sum_{k=1}^{\infty} a_k$  does not converge absolutely. Prove this theorem.

10. For  $1 \geq x \geq 0$ , and  $p \geq 1$ , show that  $(1-x)^p \geq 1-px$ . **Hint:** This can be done using the mean value theorem from calculus. Define  $f(x) \equiv (1-x)^p - 1 + px$  and show that  $f(0) = 0$  while  $f'(x) \geq 0$  for all  $x \in (0, 1)$ .
11. Using the result of Problem 10 establish Raabe's Test, an interesting variation on the ratio test. This test says the following. Suppose there exists a constant,  $C$  and a number  $p$  such that

$$\left| \frac{a_{k+1}}{a_k} \right| \leq 1 - \frac{p}{k+C}$$

for all  $k$  large enough. Then if  $p > 1$ , it follows that  $\sum_{k=1}^{\infty} a_k$  converges absolutely. **Hint:** Let  $b_k \equiv k-1+C$  and note that for all  $k$  large enough,  $b_k > 1$ . Now conclude that there exists an integer,  $k_0$  such that  $b_{k_0} > 1$  and for all  $k \geq k_0$  the given inequality above holds. Use Problem 10 to conclude that

$$\left| \frac{a_{k+1}}{a_k} \right| \leq 1 - \frac{p}{k+C} \leq \left( 1 - \frac{1}{b_{k+1}} \right)^p = \left( \frac{b_k}{b_{k+1}} \right)^p$$

showing that  $|a_k| b_k^p$  is decreasing for  $k \geq k_0$ . Thus  $|a_k| \leq C/b_k^p = C/(k-1+C)^p$ . Now use comparison theorems and the  $p$  series to obtain the conclusion of the theorem.

12. Consider the series  $\sum_{k=0}^{\infty} (-1)^k \frac{1}{\sqrt{k+1}}$ . Show this series converges and so it makes sense to write  $\left( \sum_{k=0}^{\infty} (-1)^k \frac{1}{\sqrt{k+1}} \right)^2$ . What about the Cauchy product of this series? Does it even converge? What does this mean about using algebra on infinite sums as though they were finite sums?
13. Suppose  $f$  is a nonnegative continuous decreasing function defined on  $[1, \infty)$ . Show the improper integral,  $\int_1^{\infty} f(t) dt$  and the sum  $\sum_{k=1}^{\infty} f(k)$  converge or diverge together. This is called the integral test. Use this test to verify convergence of  $\sum_{k=1}^{\infty} \frac{1}{k^\alpha}$  whenever  $\alpha > 1$  and divergence whenever  $\alpha \leq 1$ .
14. For  $p$  a positive number, determine the convergence of  $\sum_{n=2}^{\infty} \frac{1}{n(\ln(n))^p}$  for various values of  $p$ .
15. Determine the convergence of the series  $\sum_{n=1}^{\infty} \left( \sum_{k=1}^n \frac{1}{k} \right)^{-n/2}$ .
16. Verify Theorem 11.3.31 on the two series  $\sum_{k=0}^{\infty} 2^{-k}$  and  $\sum_{k=0}^{\infty} 3^{-k}$ .

## 11.5 Taylor Series

Earlier Taylor polynomials were used to approximate known functions such as  $\sin x$  and  $\ln(1+x)$ . A much more exciting idea is to use infinite series of known functions as definitions of possibly new functions.

**Definition 11.5.1** Let  $\{a_k\}_{k=0}^{\infty}$  be a sequence of numbers. The expression,

$$\sum_{k=0}^{\infty} a_k (x-a)^k \tag{11.11}$$

is called a Taylor series centered at  $a$ . This is also called a power series centered at  $a$ .

In the above definition,  $x$  is a variable. Thus you can put in various values of  $x$  and ask whether the resulting series of numbers converges. Defining,  $D$  to be the set of all values of  $x$  such that the resulting series does converge, define a new function,  $f$  defined on  $D$  as

$$f(x) \equiv \sum_{k=0}^{\infty} a_k (x-a)^k.$$

This might be a totally new function, one which has no name. Nevertheless, much can be said about such functions. The following lemma is fundamental in considering the form of  $D$  which always turns out to be an interval centered at  $a$  which may or may not contain either end point.

**Lemma 11.5.2** *Suppose  $z \in D$ . Then if  $|x-a| < |z-a|$ , then  $x \in D$  also and furthermore, the series  $\sum_{k=0}^{\infty} |a_k| |x-a|^k$  converges.*

**Proof:** Let  $1 > r = |x-a|/|z-a|$ . The  $n^{\text{th}}$  term test implies

$$\lim_{n \rightarrow \infty} |a_n| |z-a|^n = 0$$

and so for all  $n$  large enough,

$$|a_n| |z-a|^n < 1$$

so for such  $n$ ,

$$|a_n| |x-a|^n = |a_n| |z-a|^n \frac{|x-a|^n}{|z-a|^n} \leq \frac{|x-a|^n}{|z-a|^n} < r^n$$

Therefore,  $\sum_{k=0}^{\infty} |a_k| |x-a|^k$  converges by comparison with the geometric series,  $\sum r^n$ .

With this lemma, the following fundamental theorem is obtained.

**Theorem 11.5.3** *Let  $\sum_{k=0}^{\infty} a_k (x-a)^k$  be a Taylor series. Then there exists  $r \leq \infty$  such that the Taylor series converges absolutely if  $|x-a| < r$ . Furthermore, if  $|x-a| > r$ , the Taylor series diverges.*

**Proof:** Let

$$r \equiv \sup \{|y-a| : y \in D\}.$$

Then if  $|x-a| < r$ , it follows there exists  $z \in D$  such that  $|z-a| > |x-a|$  since otherwise,  $r$  wouldn't be as defined. In fact  $|x-a|$  would then be an upper bound to  $\{|y-a| : y \in D\}$ . Therefore, by the above lemma  $\sum_{k=0}^{\infty} |a_k| |x-a|^k$  converges and this proves the first part of this theorem.

Now suppose  $|x-a| > r$ . If  $\sum_{k=0}^{\infty} a_k (x-a)^k$  converges then by the above lemma,  $r$  fails to be an upper bound to  $\{|y-a| : y \in D\}$  and so the Taylor series must diverge as claimed. This proves the theorem.

From now on  $D$  will be referred to as the interval of convergence and  $r$  of the above theorem as the radius of convergence. Determining which points of  $\{x : |x-a| = r\}$  are in  $D$  requires the use of specific convergence tests and can be quite hard. However, the determination of  $r$  tends to be pretty easy.

**Example 11.5.4** *Find the interval of convergence of the Taylor series  $\sum_{n=1}^{\infty} \frac{x^n}{n}$ .*

Use Corollary 11.3.22.

$$\lim_{n \rightarrow \infty} \left( \frac{|x|^n}{n} \right)^{1/n} = \lim_{n \rightarrow \infty} \frac{|x|}{\sqrt[n]{n}} = |x|$$

because  $\lim_{n \rightarrow \infty} \sqrt[n]{n} = 1$  and so if  $|x| < 1$  the series converges. The endpoints require special attention. When  $x = 1$  the series diverges because it reduces to  $\sum_{n=1}^{\infty} \frac{1}{n}$ . At the other endpoint, however, the series converges because it reduces to  $\sum_{n=1}^{\infty} \frac{(-1)^n}{n}$  and the alternating series test applies and gives convergence.

**Example 11.5.5** Find the radius of convergence of  $\sum_{n=1}^{\infty} \frac{n^n}{n!} x^n$ .

Apply the ratio test. Taking the ratio of the absolute values of the  $(n+1)^{th}$  and the  $n^{th}$  terms

$$\frac{\frac{(n+1)^{(n+1)}}{(n+1)n!} |x|^{n+1}}{\frac{n^n}{n!} |x|^n} = (n+1)^n |x| n^{-n} = |x| \left(1 + \frac{1}{n}\right)^n \rightarrow |x| e$$

Therefore the series converges absolutely if  $|x| e < 1$  and diverges if  $|x| e > 1$ . Consequently,  $r = 1/e$ .

### 11.5.1 Operations On Power Series

It is desirable to be able to differentiate, integrate, and multiply power series. The following theorem says one can differentiate power series in the most natural way on the interval of convergence, just as you would differentiate a polynomial. This theorem may seem obvious, but it is a serious mistake to think this. You usually cannot differentiate an infinite series whose terms are functions even if the functions are themselves polynomials. The following is special and pertains to power series. It is another example of the interchange of two limits, in this case, the limit involved in taking the derivative and the limit of the sequence of finite sums.

**Theorem 11.5.6** Let  $\sum_{n=0}^{\infty} a_n (x-a)^n$  be a Taylor series having radius of convergence  $r > 0$  and let

$$f(x) \equiv \sum_{n=0}^{\infty} a_n (x-a)^n \quad (11.12)$$

for  $|x-a| < r$ . Then

$$f'(x) = \sum_{n=0}^{\infty} a_n n (x-a)^{n-1} = \sum_{n=1}^{\infty} a_n n (x-a)^{n-1} \quad (11.13)$$

and this new differentiated power series, the derived series, has radius of convergence equal to  $r$ .

**Proof:** First it will be shown that the series on the right in (11.13) has the same radius of convergence as the original series. Thus let  $|x-a| < r$  and pick  $y$  such that

$$|x-a| < |y-a| < r.$$

Then

$$\lim_{n \rightarrow \infty} |a_n| |y-a|^{n-1} = \lim_{n \rightarrow \infty} |a_n| |y-a|^n = 0$$

because

$$\sum_{n=0}^{\infty} |a_n| |y-a|^n < \infty$$

and so, for  $n$  large enough,

$$|a_n| |y-a|^{n-1} < 1.$$

Therefore, for large enough  $n$ ,

$$\begin{aligned} |a_n| n |x - a|^{n-1} &= |a_n| |y - a|^{n-1} n \left| \frac{x - a}{y - a} \right|^{n-1} \\ &\leq n \left| \frac{x - a}{y - a} \right|^{n-1} \end{aligned}$$

and

$$\sum_{n=1}^{\infty} n \left| \frac{x - a}{y - a} \right|^{n-1}$$

converges by the ratio test. By the comparison test, it follows  $\sum_{n=1}^{\infty} a_n n (x - a)^{n-1}$  converges absolutely for any  $x$  satisfying  $|x - a| < r$ . Therefore, the radius of convergence of the derived series is at least as large as that of the original series. On the other hand, if  $\sum_{n=1}^{\infty} |a_n| n |x - a|^{n-1}$  converges then by the comparison test,  $\sum_{n=1}^{\infty} |a_n| |x - a|^{n-1}$  and therefore  $\sum_{n=1}^{\infty} |a_n| |x - a|^n$  also converges which shows the radius of convergence of the derived series is no larger than that of the original series. It remains to verify the assertion about the derivative.

Let  $|x - a| < r$  and let  $r_1 < r$  be close enough to  $r$  that

$$x \in (a - r_1, a + r_1) \subseteq [a - r_1, a + r_1] \subseteq (a - r, a + r).$$

Thus, letting  $r_2 \in (r_1, r)$ ,

$$\sum_{n=0}^{\infty} |a_n| r_1^n, \sum_{n=0}^{\infty} |a_n| r_2^n < \infty \quad (11.14)$$

Letting  $y$  be close enough to  $x$ , it follows both  $x$  and  $y$  are in  $[a - r_1, a + r_1]$ . Then considering the difference quotient,

$$\begin{aligned} \frac{f(y) - f(x)}{y - x} &= \sum_{n=0}^{\infty} a_n (y - x)^{-1} [(y - a)^n - (x - a)^n] \\ &= \sum_{n=1}^{\infty} a_n n z_n^{n-1} \end{aligned} \quad (11.15)$$

where the last equation follows from the mean value theorem and  $z_n$  is some point between  $x - a$  and  $y - a$ . Therefore,

$$\begin{aligned} \frac{f(y) - f(x)}{y - x} &= \sum_{n=1}^{\infty} a_n n z_n^{n-1} = \\ &= \sum_{n=1}^{\infty} a_n n (z_n^{n-1} - (x - a)^{n-1}) + \sum_{n=1}^{\infty} a_n n (x - a)^{n-1} \\ &= \sum_{n=2}^{\infty} a_n n (n - 1) w_n^{n-2} (z_n - (x - a)) + \sum_{n=1}^{\infty} a_n n (x - a)^{n-1} \end{aligned} \quad (11.16)$$

where  $w_n$  is between  $z_n$  and  $x - a$ . Thus  $w_n$  is between  $x - a$  and  $y - a$  and so

$$w_n + a \in [a - r_1, a + r_1]$$

which implies  $|w_n| \leq r_1$ . The first sum on the right in (11.16) therefore satisfies

$$\begin{aligned} \left| \sum_{n=2}^{\infty} a_n n(n-1) w_n^{n-2} (z_n - (x-a)) \right| &\leq |y-x| \sum_{n=2}^{\infty} |a_n| n(n-1) |w_n|^{n-2} \\ &\leq |y-x| \sum_{n=2}^{\infty} |a_n| n(n-1) r_1^{n-2} \\ &= |y-x| \sum_{n=2}^{\infty} |a_n| r_2^{n-2} n(n-1) \left( \frac{r_1}{r_2} \right)^{n-2} \end{aligned}$$

Now from (11.14),  $|a_n| r_2^{n-2} < 1$  for all  $n$  large enough. Therefore, for such  $n$ ,

$$|a_n| r_2^{n-2} n(n-1) \left( \frac{r_1}{r_2} \right)^{n-2} \leq n(n-1) \left( \frac{r_1}{r_2} \right)^{n-2}$$

and the series  $\sum n(n-1) \left( \frac{r_1}{r_2} \right)^{n-2}$  converges by the ratio test. Therefore, there exists a constant,  $C$  independent of  $y$  such that

$$\sum_{n=2}^{\infty} |a_n| n(n-1) r_1^{n-2} = C < \infty$$

Consequently, from (11.16)

$$\left| \frac{f(y) - f(x)}{y - x} - \sum_{n=1}^{\infty} a_n n (x-a)^{n-1} \right| \leq C |y-x|.$$

Taking the limit as  $y \rightarrow x$  (11.13) follows. This proves the theorem.

As an immediate corollary, it is possible to characterize the coefficients of a Taylor series.

**Corollary 11.5.7** *Let  $\sum_{n=0}^{\infty} a_n (x-a)^n$  be a Taylor series with radius of convergence  $r > 0$  and let*

$$f(x) \equiv \sum_{n=0}^{\infty} a_n (x-a)^n. \quad (11.17)$$

*Then*

$$a_n = \frac{f^{(n)}(a)}{n!}. \quad (11.18)$$

**Proof:** From (11.17),  $f(a) = a_0 \equiv f^{(0)}(a)/0!$ . From Theorem 11.5.6,

$$f'(x) = \sum_{n=1}^{\infty} a_n n (x-a)^{n-1} = a_1 + \sum_{n=2}^{\infty} a_n n (x-a)^{n-1}.$$

Now let  $x = a$  and obtain that  $f'(a) = a_1 = f'(a)/1!$ . Next use Theorem 11.5.6 again to take the second derivative and obtain

$$f''(x) = 2a_2 + \sum_{n=3}^{\infty} a_n n(n-1) (x-a)^{n-2}$$

let  $x = a$  in this equation and obtain  $a_2 = f''(a)/2 = f''(a)/2!$ . Continuing this way proves the corollary.

This also shows the coefficients of a Taylor series are unique. That is, if

$$\sum_{k=0}^{\infty} a_k (x-a)^k = \sum_{k=0}^{\infty} b_k (x-a)^k$$

for all  $x$  in some interval, then  $a_k = b_k$  for all  $k$ .

**Example 11.5.8** Find the power series for  $\sin(x)$ , and  $\cos(x)$  centered at 0 and give the interval of convergence.

First consider  $f(x) = \sin(x)$ . Then  $f'(x) = \cos(x)$ ,  $f''(x) = -\sin(x)$ ,  $f'''(x) = -\cos(x)$  etc. Therefore, from Taylor's formula, Theorem 11.1.1 on Page 257,

$$f(x) = 0 + x + 0 - \frac{x^3}{3!} + 0 + \frac{x^5}{5!} + \cdots + \frac{x^{2n+1}}{(2n+1)!} + \frac{f^{(2n+2)}(\xi_n)}{(2n+2)!}$$

where  $\xi_n$  is some number between 0 and  $x$ . Furthermore, this equals either  $\pm \sin(\xi_n)$  or  $\pm \cos(\xi_n)$  and so its absolute value is no larger than 1. Thus

$$\left| \frac{f^{(2n+2)}(\xi_n)}{(2n+2)!} \right| \leq \frac{1}{(2n+2)!}.$$

By the ratio test, it follows that

$$\sum_{n=0}^{\infty} \frac{1}{(2n+2)!} < \infty$$

and so by the comparison test,

$$\sum_{n=0}^{\infty} \left| \frac{f^{(2n+2)}(\xi_n)}{(2n+2)!} \right| < \infty$$

also. Therefore, by the  $n^{\text{th}}$  term test  $\lim_{n \rightarrow \infty} \frac{f^{(2n+2)}(\xi_n)}{(2n+2)!} = 0$ . This implies

$$\sin(x) = \sum_{k=0}^n (-1)^k \frac{x^{2k+1}}{(2k+1)!} + \frac{f^{(2n+2)}(\xi_n)}{(2n+2)!}$$

and the last term converges to zero as  $n \rightarrow \infty$  for any value of  $x$  and therefore,

$$\sin(x) = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k+1}}{(2k+1)!}$$

for all  $x \in \mathbb{R}$ . By Theorem 11.5.6, you can differentiate both sides, doing the series term by term and obtain

$$\cos(x) = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k}}{(2k)!}$$

for all  $x \in \mathbb{R}$ .

**Example 11.5.9** Find the sum  $\sum_{k=1}^{\infty} k2^{-k}$ .



It may not be obvious what this sum equals but with the above theorem it is easy to find. From the formula for the sum of a geometric series,  $\frac{1}{1-t} = \sum_{k=0}^{\infty} t^k$  if  $|t| < 1$ . Differentiate both sides to obtain

$$(1-t)^{-2} = \sum_{k=1}^{\infty} k t^{k-1}$$

whenever  $|t| < 1$ . Let  $t = 1/2$ . Then

$$4 = \frac{1}{(1-(1/2))^2} = \sum_{k=1}^{\infty} k 2^{-(k-1)}$$

and so if you multiply both sides by  $2^{-1}$ ,

$$2 = \sum_{k=1}^{\infty} k 2^{-k}.$$

The following is a very important example known as the binomial series.

**Example 11.5.10** Find a Taylor series for the function  $(1+x)^\alpha$  centered at 0 valid for  $|x| < 1$ .

Use Theorem 11.5.6 to do this. First note that if  $y(x) \equiv (1+x)^\alpha$ , then  $y$  is a solution of the following initial value problem.

$$y' - \frac{\alpha}{(1+x)}y = 0, \quad y(0) = 1. \quad (11.19)$$

Next it is necessary to observe there is only one solution to this initial value problem. To see this, multiply both sides of the differential equation in (11.19) by  $(1+x)^{-\alpha}$ . When this is done one obtains

$$\frac{d}{dx} \left( (1+x)^{-\alpha} y \right) = (1+x)^{-\alpha} \left( y' - \frac{\alpha}{(1+x)}y \right) = 0. \quad (11.20)$$

Therefore, from (11.20), there must exist a constant,  $C$ , such that

$$(1+x)^{-\alpha} y = C.$$

However,  $y(0) = 1$  and so it must be that  $C = 1$ . Therefore, there is exactly one solution to the initial value problem in (11.19) and it is  $y(x) = (1+x)^\alpha$ . The strategy for finding the Taylor series of this function consists of finding a series which solves the initial value problem above. Let

$$y(x) \equiv \sum_{n=0}^{\infty} a_n x^n \quad (11.21)$$

be a solution to (11.19). Of course it is not known at this time whether such a series exists. However, the process of finding it will demonstrate its existence. From Theorem 11.5.6 and the initial value problem,

$$(1+x) \sum_{n=0}^{\infty} a_n n x^{n-1} - \sum_{n=0}^{\infty} \alpha a_n x^n = 0$$

and so

$$\sum_{n=1}^{\infty} a_n n x^{n-1} + \sum_{n=0}^{\infty} a_n (n - \alpha) x^n = 0$$

Changing the order variable of summation in the first sum,

$$\sum_{n=0}^{\infty} a_{n+1} (n+1) x^n + \sum_{n=0}^{\infty} a_n (n-\alpha) x^n = 0$$

and from Corollary 11.5.7 and the initial condition for (11.19) this requires

$$a_{n+1} = \frac{a_n (\alpha - n)}{n+1}, a_0 = 1. \quad (11.22)$$

Therefore, from (11.22) and letting  $n = 0$ ,  $a_1 = \alpha$ . Then using (11.22) again along with this information,  $a_2 = \frac{\alpha(\alpha-1)}{2}$ . Using the same process,  $a_3 = \frac{(\frac{\alpha(\alpha-1)}{2})(\alpha-2)}{3} = \frac{\alpha(\alpha-1)(\alpha-2)}{3!}$ . By now you can spot the pattern. In general,

$$a_n = \frac{\overbrace{\alpha(\alpha-1) \cdots (\alpha-n+1)}^{n \text{ of these factors}}}{n!}.$$

Therefore, our candidate for the Taylor series is

$$y(x) = \sum_{n=0}^{\infty} \frac{\alpha(\alpha-1) \cdots (\alpha-n+1)}{n!} x^n.$$

Furthermore, the above discussion shows this series solves the initial value problem on its interval of convergence. It only remains to show the radius of convergence of this series equals 1. It will then follow that this series equals  $(1+x)^\alpha$  because of uniqueness of the initial value problem. To find the radius of convergence, use the ratio test. Thus the ratio of the absolute values of  $(n+1)^{st}$  term to the absolute value of the  $n^{th}$  term is

$$\frac{\left| \frac{\alpha(\alpha-1) \cdots (\alpha-n+1)(\alpha-n)}{(n+1)n!} \right| |x|^{n+1}}{\left| \frac{\alpha(\alpha-1) \cdots (\alpha-n+1)}{n!} \right| |x|^n} = |x| \frac{|\alpha-n|}{n+1} \rightarrow |x|$$

showing that the radius of convergence is 1 since the series converges if  $|x| < 1$  and diverges if  $|x| > 1$ .

The expression,  $\frac{\alpha(\alpha-1) \cdots (\alpha-n+1)}{n!}$  is often denoted as  $\binom{\alpha}{n}$ . With this notation, the following theorem has been established.

**Theorem 11.5.11** *Let  $\alpha$  be a real number and let  $|x| < 1$ . Then*

$$(1+x)^\alpha = \sum_{n=0}^{\infty} \binom{\alpha}{n} x^n.$$

There is a very interesting issue related to the above theorem which illustrates the limitation of power series. The function  $f(x) = (1+x)^\alpha$  makes sense for all  $x > -1$  but one is only able to describe it with a power series on the interval  $(-1, 1)$ . Think about this. The above technique is a standard one for obtaining solutions of differential equations and this example illustrates a deficiency in the method. To completely understand power series, it is necessary to take a course in complex analysis. You may have noticed the prominent role played by geometric series. This is no accident. It turns out that the right way to consider Taylor series is through the use of geometric series and something called the Cauchy integral formula of complex analysis. However, these are topics for another course.

You can also integrate power series on their interval of convergence.

**Theorem 11.5.12** Let  $f(x) = \sum_{n=0}^{\infty} a_n (x-a)^n$  and suppose the interval of convergence is  $r > 0$ . Then if  $|y-a| < r$ ,

$$\int_a^y f(x) dx = \sum_{n=0}^{\infty} \int_a^y a_n (x-a)^n dx = \sum_{n=0}^{\infty} \frac{a_n (y-a)^{n+1}}{n+1}.$$

**Proof:** Define  $F(y) \equiv \int_a^y f(x) dx$  and  $G(y) \equiv \sum_{n=0}^{\infty} \frac{a_n (y-a)^{n+1}}{n+1}$ . By Theorem 11.5.6 and the Fundamental theorem of calculus,

$$G'(y) = \sum_{n=0}^{\infty} a_n (y-a)^n = f(y) = F'(y).$$

Therefore,  $G(y) - F(y) = C$  for some constant. But  $C = 0$  because  $F(a) - G(a) = 0$ . This proves the theorem.

Next consider the problem of multiplying two power series.

**Theorem 11.5.13** Let  $\sum_{n=0}^{\infty} a_n (x-a)^n$  and  $\sum_{n=0}^{\infty} b_n (x-a)^n$  be two power series having radii of convergence  $r_1$  and  $r_2$ , both positive. Then

$$\left( \sum_{n=0}^{\infty} a_n (x-a)^n \right) \left( \sum_{n=0}^{\infty} b_n (x-a)^n \right) = \sum_{n=0}^{\infty} \left( \sum_{k=0}^n a_k b_{n-k} \right) (x-a)^n$$

whenever  $|x-a| < r \equiv \min(r_1, r_2)$ .

**Proof:** By Theorem 11.5.3 both series converge absolutely if  $|x-a| < r$ . Therefore, by Theorem 11.3.31

$$\begin{aligned} & \left( \sum_{n=0}^{\infty} a_n (x-a)^n \right) \left( \sum_{n=0}^{\infty} b_n (x-a)^n \right) = \\ & \sum_{n=0}^{\infty} \sum_{k=0}^n a_k (x-a)^k b_{n-k} (x-a)^{n-k} = \sum_{n=0}^{\infty} \left( \sum_{k=0}^n a_k b_{n-k} \right) (x-a)^n. \end{aligned}$$

This proves the theorem.

The significance of this theorem in terms of applications is that it states you can multiply power series just as you would multiply polynomials and everything will be all right on the common interval of convergence.

This theorem can be used to find Taylor series which would perhaps be hard to find without it. Here is an example.

**Example 11.5.14** Find the Taylor series for  $e^x \sin x$  centered at  $x = 0$ .

Using Problems 7 - 9 on Page 260 or Example 11.5.8 on Page 280, and Example 11.3.23 on Page 268, all that is required is to multiply

$$\left( \overbrace{1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \cdots}^{e^x} \right) \left( \overbrace{x - \frac{x^3}{3!} + \frac{x^5}{5!} + \cdots}^{\sin x} \right)$$

From the above theorem the result should be

$$x + x^2 + \left( -\frac{1}{3!} + \frac{1}{2!} \right) x^3 + \cdots$$

$$= x + x^2 + \frac{1}{3}x^3 + \cdots$$

You can continue this way and get the following to a few more terms.

$$x + x^2 + \frac{1}{3}x^3 - \frac{1}{30}x^5 - \frac{1}{90}x^6 - \frac{1}{630}x^7 + \cdots$$

I don't see a pattern in these coefficients but I can go on generating them as long as I want. (In practice this tends to not be very long.) I also know the resulting power series will converge for all  $x$  because both the series for  $e^x$  and the one for  $\sin x$  converge for all  $x$ .

**Example 11.5.15** Find the Taylor series for  $\tan x$  centered at  $x = 0$ .

Lets suppose it has a Taylor series  $a_0 + a_1x + a_2x^2 + \cdots$ . Then

$$(a_0 + a_1x + a_2x^2 + \cdots) \left( \overbrace{1 - \frac{x^2}{2} + \frac{x^4}{4!} + \cdots}^{\cos x} \right) = \left( x - \frac{x^3}{3!} + \frac{x^5}{5!} + \cdots \right).$$

Using the above,  $a_0 = 0, a_1x = x$  so  $a_1 = 1, (0 - \frac{1}{2}) + a_2)x^2 = 0$  so  $a_2 = 0$ .  $(a_3 - \frac{a_1}{2})x^3 = \frac{-1}{3!}x^3$  so  $a_3 - \frac{1}{2} = -\frac{1}{6}$  so  $a_3 = \frac{1}{3}$ . Clearly one can continue in this manner. Thus the first several terms of the power series for  $\tan$  are

$$\tan x = x + \frac{1}{3}x^3 + \cdots$$

You can go on calculating these terms and find the next two yielding

$$\tan x = x + \frac{1}{3}x^3 + \frac{2}{15}x^5 + \frac{17}{315}x^7 + \cdots$$

This is a very significant technique because, as you see, there does not appear to be a very simple pattern for the coefficients of the power series for  $\tan x$ . Of course there are some issues here about whether  $\tan x$  even has a power series, but if it does, the above must be it. In fact,  $\tan(x)$  will have a power series valid on some interval centered at 0 and this becomes completely obvious when one uses methods from complex analysis but it isn't too obvious at this point. If you are interested in this issue, read the last section of the chapter. Note also that what has been accomplished is to divide the power series for  $\sin x$  by the power series for  $\cos x$  just like they were polynomials.

## 11.6 Exercises

1. Find the radius of convergence of the following.

- (a)  $\sum_{k=1}^{\infty} \left(\frac{x}{2}\right)^n$
- (b)  $\sum_{k=1}^{\infty} \sin\left(\frac{1}{n}\right) 3^n x^n$
- (c)  $\sum_{k=0}^{\infty} k! x^k$
- (d)  $\sum_{n=0}^{\infty} \frac{(3n)^n}{(3n)!} x^n$
- (e)  $\sum_{n=0}^{\infty} \frac{(2n)^n}{(2n)!} x^n$

2. Find  $\sum_{k=1}^{\infty} k2^{-k}$ .

3. Find the power series centered at 0 for the function  $1/(1+x^2)$  and give the radius of convergence.
4. Use the power series technique which was applied in Example 11.5.10 to consider the initial value problem  $y' = y, y(0) = 1$ . This yields another way to obtain the power series for  $e^x$ .
5. Use the power series technique on the initial value problem  $y' + y = 0, y(0) = 1$ . What is the solution to this initial value problem?
6. Use the power series technique to find solutions in terms of power series to the initial value problem

$$y'' + xy = 0, y(0) = 0, y'(0) = 1.$$

Tell where your solution gives a valid description of a solution for the initial value problem. **Hint:** This is a little different but you proceed the same way as in Example 11.5.10. The main difference is you have to do two differentiations of the power series instead of one.

7. Suppose the function,  $e^x$  is defined in terms of a power series,  $e^x \equiv \sum_{k=0}^{\infty} \frac{x^k}{k!}$ . Use Theorem 11.3.31 on Page 273 to show directly the usual law of exponents,

$$e^{x+y} = e^x e^y.$$

Be sure to check all the hypotheses.

8. Define the following function<sup>2</sup>:

$$f(x) \equiv \begin{cases} e^{-(1/x^2)} & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases}.$$

Show that  $f^{(k)}(x)$  exists for all  $k$  and for all  $x$ . Show also that  $f^{(k)}(0) = 0$  for all  $k \in \mathbb{N}$ . Therefore, the power series for  $f(x)$  is of the form  $\sum_{k=0}^{\infty} 0x^k$  and it converges for all values of  $x$ . However, it fails to converge to  $f(x)$  except at the single point,  $x = 0$ .

9. Let  $f_n(x) \equiv (\frac{1}{n} + x^2)^{1/2}$ . Show that for all  $x$ ,

$$||x| - f_n(x)| \leq \frac{1}{\sqrt{n}}.$$

Now show  $f'_n(0) = 0$  for all  $n$  and so  $f'_n(0) \rightarrow 0$ . However, the function,  $f(x) \equiv |x|$  has no derivative at  $x = 0$ . Thus even though  $f_n(x) \rightarrow f(x)$  for all  $x$ , you cannot say that  $f'_n(0) \rightarrow f'(0)$ .

10. Let the functions,  $f_n(x)$  be given in Problem 9 and consider

$$g_1(x) = f_1(x), g_n(x) = f_n(x) - f_{n-1}(x) \text{ if } n > 1.$$

---

<sup>2</sup>Surprisingly, this function is very important to those who use modern techniques to study differential equations. One needs to consider test functions which have the property they have infinitely many derivatives but vanish outside of some interval. The theory of complex variables can be used to show there are no examples of such functions if they have a valid power series expansion. It even becomes a little questionable whether such strange functions even exist at all. Nevertheless, they do, there are enough of them, and it is this very example which is used to show this.

Show that for all  $x$ ,

$$\sum_{k=0}^{\infty} g_k(x) = |x|$$

and that  $g'_k(0) = 0$  for all  $k$ . Therefore, you can't differentiate the series term by term and get the right answer<sup>3</sup>.

11. Find the exact value of  $\sum_{n=1}^{\infty} n^2 2^{-n}$ .
12. Use the theorem about the binomial series to give a proof of the binomial theorem

$$(a+b)^n = \sum_{k=0}^n \binom{n}{k} a^{n-k} b^k$$

whenever  $n$  is a positive integer.

13. You know  $\int_0^x \frac{1}{t+1} dt = \ln|1+x|$ . Use this and Theorem 11.5.12 to find the power series for  $\ln|1+x|$  centered at 0. Where does this power series converge? Where does it converge to the function,  $\ln|1+x|$ ?
14. You know  $\int_0^x \frac{1}{t^2+1} dt = \arctan x$ . Use this and Theorem 11.5.12 to find the power series for  $\arctan x$  centered at 0. Where does this power series converge? Where does it converge to the function,  $\arctan x$ ?
15. Find the power series for  $\sin(x^2)$  by plugging in  $x^2$  where ever there is an  $x$  in the power series for  $\sin x$ . How do you know this is the power series for  $\sin(x^2)$ ?
16. It is hard to find  $\int_0^1 e^{x^2} dx$  because you don't have a convenient antiderivative for the integrand. Replace  $e^{x^2}$  with an appropriate power series and estimate this integral.
17. Do the same as the previous problem for  $\int_0^1 \sin(x^2) dx$ .
18. Find  $\lim_{x \rightarrow 0} \frac{\tan(\sin x) - \sin(\tan x)}{x^7}$ .<sup>4</sup>

## 11.7 Some Other Theorems

First recall Theorem 11.3.31 on Page 273. For convenience, the version of this theorem which is of interest here is listed below.

**Theorem 11.7.1** Suppose  $\sum_{i=0}^{\infty} a_i$  and  $\sum_{j=0}^{\infty} b_j$  both converge absolutely. Then

$$\left( \sum_{i=0}^{\infty} a_i \right) \left( \sum_{j=0}^{\infty} b_j \right) = \sum_{n=0}^{\infty} c_n$$

where

$$c_n = \sum_{k=0}^n a_k b_{n-k}.$$

Furthermore,  $\sum_{n=0}^{\infty} c_n$  converges absolutely.

<sup>3</sup>How bad can this get? It can be much worse than this. In fact, there are functions which are continuous everywhere and differentiable nowhere. We typically don't have names for them but they are there just the same. Every such function can be written as an infinite sum of polynomials which of course have derivatives at every point. Thus it is nonsense to differentiate an infinite sum term by term without a theorem of some sort.

<sup>4</sup>This is a wonderful example. You should plug in small values of  $x$  using a calculator and see what you get using modern technology.

**Proof:** It only remains to verify the last series converges absolutely. By Theorem 11.3.28 on Page 271 and letting  $p_{nk}$  be as defined there,

$$\begin{aligned}
 \sum_{n=0}^{\infty} |c_n| &= \sum_{n=0}^{\infty} \left| \sum_{k=0}^n a_k b_{n-k} \right| \\
 &\leq \sum_{n=0}^{\infty} \sum_{k=0}^n |a_k| |b_{n-k}| = \sum_{n=0}^{\infty} \sum_{k=0}^{\infty} p_{nk} |a_k| |b_{n-k}| \\
 &= \sum_{k=0}^{\infty} \sum_{n=0}^{\infty} p_{nk} |a_k| |b_{n-k}| = \sum_{k=0}^{\infty} \sum_{n=k}^{\infty} |a_k| |b_{n-k}| \\
 &= \sum_{k=0}^{\infty} |a_k| \sum_{n=0}^{\infty} |b_n| < \infty.
 \end{aligned}$$

This proves the theorem.

The theorem is about multiplying two series. What if you wanted to consider

$$\left( \sum_{n=0}^{\infty} a_n \right)^p$$

where  $p$  is a positive integer maybe larger than 2? Is there a similar theorem to the above?

**Definition 11.7.2** Define

$$\sum_{k_1 + \cdots + k_p = m} a_{k_1} a_{k_2} \cdots a_{k_p}$$

as follows. Consider all ordered lists of nonnegative integers  $k_1, \dots, k_p$  which have the property that  $\sum_{i=1}^p k_i = m$ . For each such list of integers, form the product,  $a_{k_1} a_{k_2} \cdots a_{k_p}$  and then add all these products.

Note that

$$\sum_{k=0}^n a_k a_{n-k} = \sum_{k_1 + k_2 = n} a_{k_1} a_{k_2}$$

Therefore, from the above theorem, if  $\sum a_i$  converges absolutely, it follows

$$\left( \sum_{i=0}^{\infty} a_i \right)^2 = \sum_{n=0}^{\infty} \left( \sum_{k_1 + k_2 = n} a_{k_1} a_{k_2} \right).$$

It turns out a similar theorem holds for replacing 2 with  $p$ .

**Theorem 11.7.3** Suppose  $\sum_{n=0}^{\infty} a_n$  converges absolutely. Then

$$\left( \sum_{n=0}^{\infty} a_n \right)^p = \sum_{m=0}^{\infty} c_{mp}$$

where

$$c_{mp} \equiv \sum_{k_1 + \cdots + k_p = m} a_{k_1} \cdots a_{k_p}.$$

**Proof:** First note this is obviously true if  $p = 1$  and is also true if  $p = 2$  from the above theorem. Now suppose this is true for  $p$  and consider  $(\sum_{n=0}^{\infty} a_n)^{p+1}$ . By the induction hypothesis and the above theorem on the Cauchy product,

$$\begin{aligned}
 \left(\sum_{n=0}^{\infty} a_n\right)^{p+1} &= \left(\sum_{n=0}^{\infty} a_n\right)^p \left(\sum_{n=0}^{\infty} a_n\right) \\
 &= \left(\sum_{m=0}^{\infty} c_{mp}\right) \left(\sum_{n=0}^{\infty} a_n\right) \\
 &= \sum_{n=0}^{\infty} \left(\sum_{k=0}^n c_{kp} a_{n-k}\right) \\
 &= \sum_{n=0}^{\infty} \sum_{k=0}^n \sum_{k_1+\dots+k_p=k} a_{k_1} \cdots a_{k_p} a_{n-k} \\
 &= \sum_{n=0}^{\infty} \sum_{k_1+\dots+k_{p+1}=n} a_{k_1} \cdots a_{k_{p+1}}
 \end{aligned}$$

and this proves the theorem.

This theorem implies the following corollary for power series.

**Corollary 11.7.4** *Let*

$$\sum_{n=0}^{\infty} a_n (x-a)^n$$

*be a power series having radius of convergence,  $r > 0$ . Then if  $|x-a| < r$ ,*

$$\left(\sum_{n=0}^{\infty} a_n (x-a)^n\right)^p = \sum_{n=0}^{\infty} b_{np} (x-a)^n$$

*where*

$$b_{np} \equiv \sum_{k_1+\dots+k_p=n} a_{k_1} \cdots a_{k_p}.$$

**Proof:** Since  $|x-a| < r$ , the series,  $\sum_{n=0}^{\infty} a_n (x-a)^n$ , converges absolutely. Therefore, the above theorem applies and

$$\begin{aligned}
 \left(\sum_{n=0}^{\infty} a_n (x-a)^n\right)^p &= \\
 \sum_{n=0}^{\infty} \left(\sum_{k_1+\dots+k_p=n} a_{k_1} (x-a)^{k_1} \cdots a_{k_p} (x-a)^{k_p}\right) &= \\
 \sum_{n=0}^{\infty} \left(\sum_{k_1+\dots+k_p=n} a_{k_1} \cdots a_{k_p}\right) (x-a)^n.
 \end{aligned}$$

With this theorem it is possible to consider the question raised in Example 11.5.15 on Page 284 about the existence of the power series for  $\tan x$ . This question is clearly included in the more general question of when

$$\left(\sum_{n=0}^{\infty} a_n (x-a)^n\right)^{-1}$$



has a power series.

**Lemma 11.7.5** *Let  $f(x) = \sum_{n=0}^{\infty} a_n (x-a)^n$ , a power series having radius of convergence  $r > 0$ . Suppose also that  $f(a) = 1$ . Then there exists  $r_1 > 0$  and  $\{b_n\}$  such that for all  $|x-a| < r_1$ ,*

$$\frac{1}{f(x)} = \sum_{n=0}^{\infty} b_n (x-a)^n.$$

**Proof:** By continuity, there exists  $r_1 > 0$  such that if  $|x-a| < r_1$ , then

$$\sum_{n=1}^{\infty} |a_n| |x-a|^n < 1.$$

Now pick such an  $x$ . Then

$$\begin{aligned} \frac{1}{f(x)} &= \frac{1}{1 + \sum_{n=1}^{\infty} a_n (x-a)^n} \\ &= \frac{1}{1 + \sum_{n=0}^{\infty} c_n (x-a)^n} \end{aligned}$$

where  $c_n = a_n$  if  $n > 0$  and  $c_0 = 0$ . Then

$$\left| \sum_{n=1}^{\infty} a_n (x-a)^n \right| \leq \sum_{n=1}^{\infty} |a_n| |x-a|^n < 1 \quad (11.23)$$

and so from the formula for the sum of a geometric series,

$$\frac{1}{f(x)} = \sum_{p=0}^{\infty} \left( \sum_{n=0}^{\infty} c_n (x-a)^n \right)^p.$$

By Corollary 11.7.4, this equals

$$\sum_{p=0}^{\infty} \sum_{n=0}^{\infty} b_{np} (x-a)^n \quad (11.24)$$

where

$$b_{np} = \sum_{k_1 + \dots + k_p = n} c_{k_1} \cdots c_{k_p}.$$

Thus  $|b_{np}| \leq \sum_{k_1 + \dots + k_p = n} |c_{k_1}| \cdots |c_{k_p}| \equiv B_{np}$  and so by Theorem 11.7.3,

$$\begin{aligned} \sum_{p=0}^{\infty} \sum_{n=0}^{\infty} |b_{np}| |x-a|^n &\leq \sum_{p=0}^{\infty} \sum_{n=0}^{\infty} B_{np} |x-a|^n \\ &= \sum_{p=0}^{\infty} \left( \sum_{n=0}^{\infty} |c_n| |x-a|^n \right)^p < \infty \end{aligned}$$

by (11.23) and the formula for the sum of a geometric series. Since the series of (11.24) converges absolutely, Theorem 11.3.28 on Page 271 implies the series in (11.24) equals

$$\sum_{n=0}^{\infty} \left( \sum_{p=0}^{\infty} b_{np} \right) (x-a)^n$$

and so, letting  $\sum_{p=0}^{\infty} b_{np} \equiv b_n$ , this proves the lemma.

With this lemma, the following theorem is easy to obtain.

**Theorem 11.7.6** Let  $f(x) = \sum_{n=0}^{\infty} a_n (x-a)^n$ , a power series having radius of convergence  $r > 0$ . Suppose also that  $f(a) \neq 0$ . Then there exists  $r_1 > 0$  and  $\{b_n\}$  such that for all  $|x-a| < r_1$ ,

$$\frac{1}{f(x)} = \sum_{n=0}^{\infty} b_n (x-a)^n.$$

**Proof:** Let  $g(x) \equiv f(x)/f(a)$  so that  $g(x)$  satisfies the conditions of the above lemma. Then by that lemma, there exists  $r_1 > 0$  and a sequence,  $\{b_n\}$  such that

$$\frac{f(a)}{f(x)} = \sum_{n=0}^{\infty} b_n (x-a)^n$$

for all  $|x-a| < r_1$ . Then

$$\frac{1}{f(x)} = \sum_{n=0}^{\infty} \tilde{b}_n (x-a)^n$$

where  $\tilde{b}_n = b_n/f(a)$ . This proves the theorem.

There is a very interesting question related to  $r_1$  in this theorem. One might think that if  $|x-a| < r$ , the radius of convergence of  $f(x)$  and if  $f(x) \neq 0$  it should be possible to write  $1/f(x)$  as a power series centered at  $a$ . Unfortunately this is not true. Consider  $f(x) = 1+x^2$ . In this case  $r = \infty$  but the power series for  $1/f(x)$  converges only if  $|x| < 1$ . What happens is this,  $1/f(x)$  will have a power series that will converge for  $|x-a| < r_1$  where  $r_1$  is the distance between  $a$  and the nearest singularity or zero of  $f(x)$  in the complex plane. In the case of  $f(x) = 1+x^2$  this function has a zero at  $x = \pm i$ . This is just another instance of why the natural setting for the study of power series is the complex plane. To read more on power series, you should see the book by Apostol [1] or any text on complex variable.

**Part III**

**Vector Valued Functions**



$\mathbb{R}^n$

The notation,  $\mathbb{R}^n$  refers to the collection of ordered lists of  $n$  real numbers. More precisely, consider the following definition.

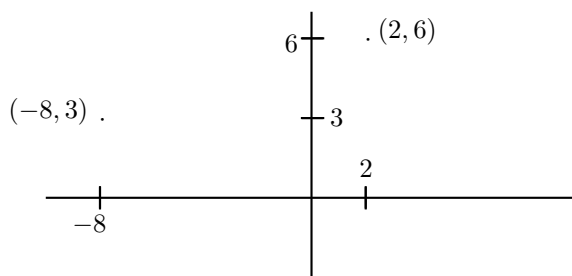
**Definition 12.0.7** Define  $\mathbb{R}^n \equiv \{(x_1, \dots, x_n) : x_j \in \mathbb{R} \text{ for } j = 1, \dots, n\}$ .  $(x_1, \dots, x_n) = (y_1, \dots, y_n)$  if and only if for all  $j = 1, \dots, n$ ,  $x_j = y_j$ . When  $(x_1, \dots, x_n) \in \mathbb{R}^n$ , it is conventional to denote  $(x_1, \dots, x_n)$  by the single bold face letter,  $\mathbf{x}$ . The numbers,  $x_j$  are called the coordinates. The set

$$\{(0, \dots, 0, t, 0, \dots, 0) : t \in \mathbb{R}\}$$

for  $t$  in the  $i^{\text{th}}$  slot is called the  $i^{\text{th}}$  coordinate axis. The point  $\mathbf{0} \equiv (0, \dots, 0)$  is called the origin.

Thus  $(1, 2, 4) \in \mathbb{R}^3$  and  $(2, 1, 4) \in \mathbb{R}^3$  but  $(1, 2, 4) \neq (2, 1, 4)$  because, even though the same numbers are involved, they don't match up. In particular, the first entries are not equal.

Why would anyone be interested in such a thing? First consider the case when  $n = 1$ . Then from the definition,  $\mathbb{R}^1 = \mathbb{R}$ . Recall that  $\mathbb{R}$  is identified with the points of a line. Look at the number line again. Observe that this amounts to identifying a point on this line with a real number. In other words a real number determines where you are on this line. Now suppose  $n = 2$  and consider two lines which intersect each other at right angles as shown in the following picture.



Notice how you can identify a point shown in the plane with the ordered pair,  $(2, 6)$ . You go to the right a distance of 2 and then up a distance of 6. Similarly, you can identify another point in the plane with the ordered pair  $(-8, 3)$ . Go to the left a distance of 8 and then up a distance of 3. The reason you go to the left is that there is a  $-$  sign on the eight. From this reasoning, every ordered pair determines a unique point in the plane. Conversely, taking a point in the plane, you could draw two lines through the point, one vertical and the

other horizontal and determine unique points,  $x_1$  on the horizontal line in the above picture and  $x_2$  on the vertical line in the above picture, such that the point of interest is identified with the ordered pair,  $(x_1, x_2)$ . In short, points in the plane can be identified with ordered pairs similar to the way that points on the real line are identified with real numbers. Now suppose  $n = 3$ . As just explained, the first two coordinates determine a point in a plane. Letting the third component determine how far up or down you go, depending on whether this number is positive or negative, this determines a point in space. Thus,  $(1, 4, -5)$  would mean to determine the point in the plane that goes with  $(1, 4)$  and then to go below this plane a distance of 5 to obtain a unique point in space. You see that the ordered triples correspond to points in space just as the ordered pairs correspond to points in a plane and single real numbers correspond to points on a line.

You can't stop here and say that you are only interested in  $n \leq 3$ . What if you were interested in the motion of two objects? You would need three coordinates to describe where the first object is and you would need another three coordinates to describe where the other object is located. Therefore, you would need to be considering  $\mathbb{R}^6$ . If the two objects moved around, you would need a time coordinate as well. As another example, consider a hot object which is cooling and suppose you want the temperature of this object. How many coordinates would be needed? You would need one for the temperature, three for the position of the point in the object and one more for the time. Thus you would need to be considering  $\mathbb{R}^5$ . Many other examples can be given. Sometimes  $n$  is very large. This is often the case in applications to business when they are trying to maximize profit subject to constraints. It also occurs in numerical analysis when people try to solve hard problems on a computer.

There are other ways to identify points in space with three numbers but the one presented is the most basic. In this case, the coordinates are known as Cartesian coordinates after Descartes<sup>1</sup> who invented this idea in the first half of the seventeenth century. I will often not bother to draw a distinction between the point in  $n$  dimensional space and its Cartesian coordinates.

## 12.1 Algebra in $\mathbb{R}^n$

There are two algebraic operations done with elements of  $\mathbb{R}^n$ . One is addition and the other is multiplication by numbers, called scalars.

**Definition 12.1.1** *If  $\mathbf{x} \in \mathbb{R}^n$  and  $a$  is a number, also called a scalar. Then  $a\mathbf{x} \in \mathbb{R}^n$  is defined by*

$$a\mathbf{x} = a(x_1, \dots, x_n) \equiv (ax_1, \dots, ax_n). \quad (12.1)$$

*This is known as scalar multiplication. If  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  then  $\mathbf{x} + \mathbf{y} \in \mathbb{R}^n$  and is defined by*

$$\begin{aligned} \mathbf{x} + \mathbf{y} &= (x_1, \dots, x_n) + (y_1, \dots, y_n) \\ &\equiv (x_1 + y_1, \dots, x_n + y_n) \end{aligned} \quad (12.2)$$

With this definition, the algebraic properties satisfy the conclusions of the following theorem.

**Theorem 12.1.2** *For  $\mathbf{v}, \mathbf{w} \in \mathbb{R}^n$  and  $\alpha, \beta$  scalars, (real numbers), the following hold.*

$$\mathbf{v} + \mathbf{w} = \mathbf{w} + \mathbf{v}, \quad (12.3)$$

---

<sup>1</sup>René Descartes 1596-1650 is often credited with inventing analytic geometry although it seems the ideas were actually known much earlier. He was interested in many different subjects, physiology, chemistry, and physics being some of them. He also wrote a large book in which he tried to explain the book of Genesis scientifically. Descartes ended up dying in Sweden.

the commutative law of addition,

$$(\mathbf{v} + \mathbf{w}) + \mathbf{z} = \mathbf{v} + (\mathbf{w} + \mathbf{z}), \quad (12.4)$$

the associative law for addition,

$$\mathbf{v} + \mathbf{0} = \mathbf{v}, \quad (12.5)$$

the existence of an additive identity,

$$\mathbf{v} + (-\mathbf{v}) = \mathbf{0}, \quad (12.6)$$

the existence of an additive inverse, Also

$$\alpha(\mathbf{v} + \mathbf{w}) = \alpha\mathbf{v} + \alpha\mathbf{w}, \quad (12.7)$$

$$(\alpha + \beta)\mathbf{v} = \alpha\mathbf{v} + \beta\mathbf{v}, \quad (12.8)$$

$$\alpha(\beta\mathbf{v}) = \alpha\beta(\mathbf{v}), \quad (12.9)$$

$$1\mathbf{v} = \mathbf{v}. \quad (12.10)$$

In the above  $\mathbf{0} = (0, \dots, 0)$ .

You should verify these properties all hold. For example, consider (12.7)

$$\begin{aligned} \alpha(\mathbf{v} + \mathbf{w}) &= \alpha(v_1 + w_1, \dots, v_n + w_n) \\ &= (\alpha(v_1 + w_1), \dots, \alpha(v_n + w_n)) \\ &= (\alpha v_1 + \alpha w_1, \dots, \alpha v_n + \alpha w_n) \\ &= (\alpha v_1, \dots, \alpha v_n) + (\alpha w_1, \dots, \alpha w_n) \\ &= \alpha\mathbf{v} + \alpha\mathbf{w}. \end{aligned}$$

As usual subtraction is defined as  $\mathbf{x} - \mathbf{y} \equiv \mathbf{x} + (-\mathbf{y})$ .

## 12.2 Exercises

1. Verify all the properties (12.3)-(12.10).
2. Compute  $5(1, 2, 3, -2) + 6(2, 1, -2, 7)$ .
3. Draw a picture of the points in  $\mathbb{R}^2$  which are determined by the following ordered pairs.
  - (a)  $(1, 2)$
  - (b)  $(-2, -2)$
  - (c)  $(-2, 3)$
  - (d)  $(2, -5)$
4. Does it make sense to write  $(1, 2) + (2, 3, 1)$ ? Explain.
5. Draw a picture of the points in  $\mathbb{R}^3$  which are determined by the following ordered triples.
  - (a)  $(1, 2, 0)$
  - (b)  $(-2, -2, 1)$
  - (c)  $(-2, 3, -2)$

### 12.3 Distance in $\mathbb{R}^n$

How is distance between two points in  $\mathbb{R}^n$  defined?

**Definition 12.3.1** Let  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_n)$  be two points in  $\mathbb{R}^n$ . Then  $|\mathbf{x} - \mathbf{y}|$  indicates the distance between these points and is defined as

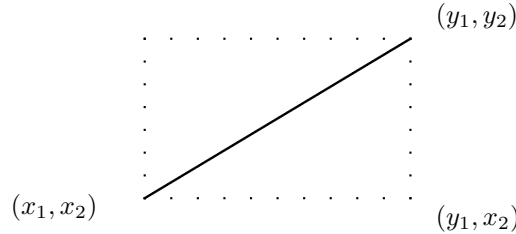
$$\text{distance between } \mathbf{x} \text{ and } \mathbf{y} \equiv |\mathbf{x} - \mathbf{y}| \equiv \left( \sum_{k=1}^n |x_k - y_k|^2 \right)^{1/2}.$$

This is called the distance formula. The symbol,  $B(\mathbf{a}, r)$  is defined by

$$B(\mathbf{a}, r) \equiv \{\mathbf{x} \in \mathbb{R}^n : |\mathbf{x} - \mathbf{a}| < r\}.$$

This is called an open ball of radius  $r$  centered at  $\mathbf{a}$ . It gives all the points in  $\mathbb{R}^n$  which are closer to  $\mathbf{a}$  than  $r$ .

First of all note this is a generalization of the notion of distance in  $\mathbb{R}$ . There the distance between two points,  $x$  and  $y$  was given by the absolute value of their difference. Thus  $|x - y|$  is equal to the distance between these two points on  $\mathbb{R}$ . Now  $|x - y| = \left( (x - y)^2 \right)^{1/2}$  where the square root is always the positive square root. Thus it is the same formula as the above definition except there is only one term in the sum. Geometrically, this is the right way to define distance which is seen from the Pythagorean theorem. Consider the following picture in the case that  $n = 2$ .



There are two points in the plane whose Cartesian coordinates are  $(x_1, x_2)$  and  $(y_1, y_2)$  respectively. Then the solid line joining these two points is the hypotenuse of a right triangle which is half of the rectangle shown in dotted lines. What is its length? Note the lengths of the sides of this triangle are  $|y_1 - x_1|$  and  $|y_2 - x_2|$ . Therefore, the Pythagorean theorem implies the length of the hypotenuse equals

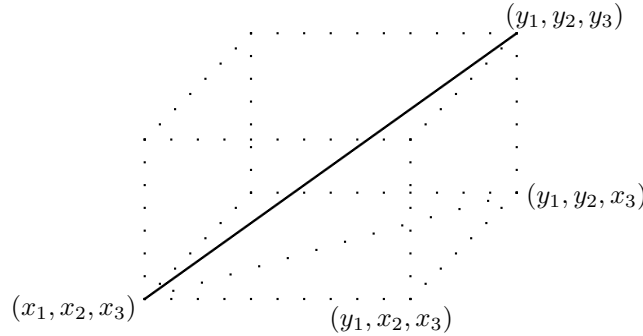
$$\left( |y_1 - x_1|^2 + |y_2 - x_2|^2 \right)^{1/2} = \left( (y_1 - x_1)^2 + (y_2 - x_2)^2 \right)^{1/2}$$

which is just the formula for the distance given above.

Now suppose  $n = 3$  and let  $(x_1, x_2, x_3)$  and  $(y_1, y_2, y_3)$  be two points in  $\mathbb{R}^3$ . Consider the following picture in which one of the solid lines joins the two points and a dotted line joins



the points  $(x_1, x_2, x_3)$  and  $(y_1, y_2, x_3)$ .



By the Pythagorean theorem, the length of the dotted line joining  $(x_1, x_2, x_3)$  and  $(y_1, y_2, x_3)$  equals

$$\left((y_1 - x_1)^2 + (y_2 - x_2)^2\right)^{1/2}$$

while the length of the line joining  $(y_1, y_2, x_3)$  to  $(y_1, y_2, y_3)$  is just  $|y_3 - x_3|$ . Therefore, by the Pythagorean theorem again, the length of the line joining the points  $(x_1, x_2, x_3)$  and  $(y_1, y_2, y_3)$  equals

$$\begin{aligned} & \left\{ \left[ \left( (y_1 - x_1)^2 + (y_2 - x_2)^2 \right)^{1/2} \right]^2 + (y_3 - x_3)^2 \right\}^{1/2} \\ &= \left( (y_1 - x_1)^2 + (y_2 - x_2)^2 + (y_3 - x_3)^2 \right)^{1/2}, \end{aligned}$$

which is again just the distance formula above.

This completes the argument that the above definition is reasonable. Of course you cannot continue drawing pictures in ever higher dimensions but there is not problem with the formula for distance in any number of dimensions. Here is an example.

**Example 12.3.2** Find the distance between the points in  $\mathbb{R}^4$ ,  $\mathbf{a} = (1, 2, -4, 6)$  and  $\mathbf{b} = (2, 3, -1, 0)$

Use the distance formula and write

$$|\mathbf{a} - \mathbf{b}|^2 = (1 - 2)^2 + (2 - 3)^2 + (-4 - (-1))^2 + (6 - 0)^2 = 47$$

Therefore,  $|\mathbf{a} - \mathbf{b}| = \sqrt{47}$ .

All this amounts to defining the distance between two points as the length of a straight line joining these two points. However, there is nothing sacred about using straight lines. One could define the distance to be the length of some other sort of line joining these points. It won't be done in this book but sometimes this sort of thing is done.

Another convention which is usually followed, especially in  $\mathbb{R}^2$  and  $\mathbb{R}^3$  is to denote the first component of a point in  $\mathbb{R}^2$  by  $x$  and the second component by  $y$ . In  $\mathbb{R}^3$  it is customary to denote the first and second components as just described while the third component is called  $z$ .

**Example 12.3.3** Describe the points which are at the same distance between  $(1, 2, 3)$  and  $(0, 1, 2)$ .

Let  $(x, y, z)$  be such a point. Then

$$\sqrt{(x-1)^2 + (y-2)^2 + (z-3)^2} = \sqrt{x^2 + (y-1)^2 + (z-2)^2}.$$

Squaring both sides

$$(x-1)^2 + (y-2)^2 + (z-3)^2 = x^2 + (y-1)^2 + (z-2)^2$$

and so

$$x^2 - 2x + 14 + y^2 - 4y + z^2 - 6z = x^2 + y^2 - 2y + 5 + z^2 - 4z$$

which implies

$$-2x + 14 - 4y - 6z = -2y + 5 - 4z$$

and so

$$2x + 2y + 2z = -9. \quad (12.11)$$

Since these steps are reversible, the set of points which is at the same distance from the two given points consists of the points,  $(x, y, z)$  such that (12.11) holds.

The following lemma is fundamental. It is a form of the Cauchy Schwarz inequality.

**Lemma 12.3.4** *Let  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_n)$  be two points in  $\mathbb{R}^n$ . Then*

$$\left| \sum_{i=1}^n x_i y_i \right| \leq |\mathbf{x}| |\mathbf{y}|. \quad (12.12)$$

**Proof:** Let  $\theta$  be either 1 or  $-1$  such that

$$\theta \sum_{i=1}^n x_i y_i = \sum_{i=1}^n x_i (\theta y_i) = \left| \sum_{i=1}^n x_i y_i \right|$$

and consider  $p(t) \equiv \sum_{i=1}^n (x_i + t\theta y_i)^2$ . Then for all  $t \in \mathbb{R}$ ,

$$\begin{aligned} 0 &\leq p(t) = \sum_{i=1}^n x_i^2 + 2t \sum_{i=1}^n x_i \theta y_i + t^2 \sum_{i=1}^n y_i^2 \\ &= |\mathbf{x}|^2 + 2t \sum_{i=1}^n x_i \theta y_i + t^2 |\mathbf{y}|^2 \end{aligned}$$

If  $|\mathbf{y}| = 0$  then (12.12) is obviously true because both sides equal zero. Therefore, assume  $|\mathbf{y}| \neq 0$  and then  $p(t)$  is a polynomial of degree two whose graph opens up. Therefore, it either has no zeroes, two zeroes or one repeated zero. If it has two zeroes, the above inequality must be violated because in this case the graph must dip below the  $x$  axis. Therefore, it either has no zeroes or exactly one. From the quadratic formula (see Problem 18 on Page 44) this happens exactly when

$$4 \left( \sum_{i=1}^n x_i \theta y_i \right)^2 - 4 |\mathbf{x}|^2 |\mathbf{y}|^2 \leq 0$$

and so

$$\sum_{i=1}^n x_i \theta y_i = \left| \sum_{i=1}^n x_i y_i \right| \leq |\mathbf{x}| |\mathbf{y}|$$

as claimed. This proves the inequality. In the next chapter a different proof is given.

There are certain properties of the distance which are obvious. Two of them which follow directly from the definition are

$$|\mathbf{x} - \mathbf{y}| = |\mathbf{y} - \mathbf{x}|,$$

$$|\mathbf{x} - \mathbf{y}| \geq 0 \text{ and equals } 0 \text{ only if } \mathbf{y} = \mathbf{x}.$$

The third fundamental property of distance is known as the triangle inequality. Recall that in any triangle the sum of the lengths of two sides is always at least as large as the third side. The following corollary is equivalent to this simple statement and this will be more clear in the next chapter. The above lemma makes possible a completely algebraic proof of this inequality.

**Corollary 12.3.5** *Let  $\mathbf{x}, \mathbf{y}$  be points of  $\mathbb{R}^n$ . Then*

$$|\mathbf{x} + \mathbf{y}| \leq |\mathbf{x}| + |\mathbf{y}|.$$

**Proof:** Using the above lemma,

$$\begin{aligned} |\mathbf{x} + \mathbf{y}|^2 &\equiv \sum_{i=1}^n (x_i + y_i)^2 \\ &= \sum_{i=1}^n x_i^2 + 2 \sum_{i=1}^n x_i y_i + \sum_{i=1}^n y_i^2 \\ &\leq |\mathbf{x}|^2 + 2|\mathbf{x}||\mathbf{y}| + |\mathbf{y}|^2 \\ &= (|\mathbf{x}| + |\mathbf{y}|)^2 \end{aligned}$$

and so upon taking square roots of both sides,

$$|\mathbf{x} + \mathbf{y}| \leq |\mathbf{x}| + |\mathbf{y}|$$

and this proves the corollary.

## 12.4 Exercises

1. You are given two points in  $\mathbb{R}^3$ ,  $(4, 5, -4)$  and  $(2, 3, 0)$ . Show the distance from the point,  $(3, 4, -2)$  to the first of these points is the same as the distance from this point to the second of the original pair of points. Note that  $3 = \frac{4+2}{2}$ ,  $4 = \frac{5+3}{2}$ . Obtain a theorem which will be valid for general pairs of points,  $(x, y, z)$  and  $(x_1, y_1, z_1)$  and prove your theorem using the distance formula.
2. A sphere is the set of all points which are at a given distance from a single given point. Find an equation for the sphere which is the set of all points that are at a distance of 4 from the point  $(1, 2, 3)$  in  $\mathbb{R}^3$ .
3. A sphere centered at the point  $(x_0, y_0, z_0) \in \mathbb{R}^3$  having radius  $r$  consists of all points,  $(x, y, z)$  whose distance to  $(x_0, y_0, z_0)$  equals  $r$ . Write an equation for this sphere in  $\mathbb{R}^3$ .
4. Suppose the distance between  $(x, y)$  and  $(x', y')$  were defined to equal the larger of the two numbers  $|x - x'|$  and  $|y - y'|$ . Draw a picture of the sphere centered at the point,  $(0, 0)$  if this notion of distance is used.

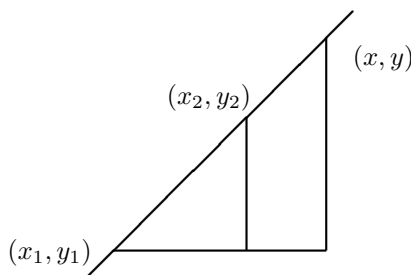
5. Repeat the same problem except this time let the distance between the two points be  $|x - x'| + |y - y'|$ .
6. If  $(x_1, y_1, z_1)$  and  $(x_2, y_2, z_2)$  are two points such that  $|(x_i, y_i, z_i)| = 1$  for  $i = 1, 2$ , show that in terms of the usual distance,  $\left| \left( \frac{x_1+x_2}{2}, \frac{y_1+y_2}{2}, \frac{z_1+z_2}{2} \right) \right| < 1$ . What would happen if you used the way of measuring distance given in Problem 4 ( $|(x, y, z)| = \text{maximum of } |z|, |x|, |y|$ )?
7. Give a simple description using the distance formula of the set of points which are at an equal distance between the two points  $(x_1, y_1, z_1)$  and  $(x_2, y_2, z_2)$ .

## 12.5 Lines in $\mathbb{R}^n$

To begin with consider the case  $n = 1, 2$ . In the case where  $n = 1$ , the only line is just  $\mathbb{R}^1 = \mathbb{R}$ . Therefore, if  $x_1$  and  $x_2$  are two different points in  $\mathbb{R}$ , consider

$$x = x_1 + t(x_2 - x_1)$$

where  $t \in \mathbb{R}$  and the totality of all such points will give  $\mathbb{R}$ . You see that you can always solve the above equation for  $t$ , showing that every point on  $\mathbb{R}$  is of this form. Now consider the plane. Does a similar formula hold? Let  $(x_1, y_1)$  and  $(x_2, y_2)$  be two different points in  $\mathbb{R}^2$  which are contained in a line,  $l$ . Suppose that  $x_1 \neq x_2$ . Then if  $(x, y)$  is an arbitrary point on  $l$ ,



Now by similar triangles,

$$m \equiv \frac{y_2 - y_1}{x_2 - x_1} = \frac{y - y_1}{x - x_1}$$

and so the point slope form of the line,  $l$ , is given as

$$y - y_1 = m(x - x_1).$$

If  $t$  is defined by

$$x = x_1 + t(x_2 - x_1),$$

you obtain this equation along with

$$\begin{aligned} y &= y_1 + mt(x_2 - x_1) \\ &= y_1 + t(y_2 - y_1). \end{aligned}$$

Therefore,

$$(x, y) = (x_1, y_1) + t(x_2 - x_1, y_2 - y_1).$$

If  $x_1 = x_2$ , then in place of the point slope form above,  $x = x_1$ . Since the two given points are different,  $y_1 \neq y_2$  and so you still obtain the above formula for the line. Because of this, the following is the definition of a line in  $\mathbb{R}^n$ .

**Definition 12.5.1** A line in  $\mathbb{R}^n$  containing the two different points,  $\mathbf{x}^1$  and  $\mathbf{x}^2$  is the collection of points of the form

$$\mathbf{x} = \mathbf{x}^1 + t(\mathbf{x}^2 - \mathbf{x}^1)$$

where  $t \in \mathbb{R}$ . This is known as a parametric equation and the variable  $t$  is called the parameter.

Often  $t$  denotes time in applications to Physics. Note this definition agrees with the usual notion of a line in two dimensions and so this is consistent with earlier concepts.

**Lemma 12.5.2** Let  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$  with  $\mathbf{a} \neq \mathbf{0}$ . Then  $\mathbf{x} = t\mathbf{a} + \mathbf{b}$ ,  $t \in \mathbb{R}$ , is a line.

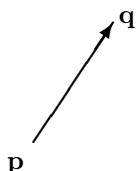
**Proof:** Let  $\mathbf{x}^1 = \mathbf{b}$  and let  $\mathbf{x}^2 - \mathbf{x}^1 = \mathbf{a}$  so that  $\mathbf{x}^2 \neq \mathbf{x}^1$ . Then  $t\mathbf{a} + \mathbf{b} = \mathbf{x}^1 + t(\mathbf{x}^2 - \mathbf{x}^1)$  and so  $\mathbf{x} = t\mathbf{a} + \mathbf{b}$  is a line containing the two different points,  $\mathbf{x}^1$  and  $\mathbf{x}^2$ . This proves the lemma.

**Definition 12.5.3** Let  $\mathbf{p}$  and  $\mathbf{q}$  be two points in  $\mathbb{R}^n$ ,  $\mathbf{p} \neq \mathbf{q}$ . The directed line segment from  $\mathbf{p}$  to  $\mathbf{q}$ , denoted by  $\overrightarrow{\mathbf{pq}}$ , is defined to be the collection of points,

$$\mathbf{x} = \mathbf{p} + t(\mathbf{q} - \mathbf{p}), t \in [0, 1]$$

with the direction corresponding to increasing  $t$ .

Think of  $\overrightarrow{\mathbf{pq}}$  as an arrow whose point is on  $\mathbf{q}$  and whose base is at  $\mathbf{p}$  as shown in the following picture.



This line segment is a part of a line from the above Definition.

**Example 12.5.4** Find a parametric equation for the line through the points  $(1, 2, 0)$  and  $(2, -4, 6)$ .

Use the definition of a line given above to write

$$(x, y, z) = (1, 2, 0) + t(1, -6, 6), t \in \mathbb{R}.$$

The reason for the word, “a”, rather than the word, “the” is there are infinitely many different parametric equations for the same line. To see this replace  $t$  with  $3s$ . Then you obtain a parametric equation for the same line because the same set of points are obtained. The difference is they are obtained from different values of the parameter. What happens is this: The line is a set of points but the parametric description gives more information than that. It tells us how the set of points are obtained. Obviously, there are many ways to trace out a given set of points and each of these ways corresponds to a different parametric equation for the line.

## 12.6 Exercises

1. Suppose you are given two points,  $(-a, 0)$  and  $(a, 0)$  in  $\mathbb{R}^2$  and a number,  $r > 2a$ . The set of points described by

$$\{(x, y) \in \mathbb{R}^2 : |(x, y) - (-a, 0)| + |(x, y) - (a, 0)| = r\}$$

is known as an ellipse. The two given points are known as the focus points of the ellipse. Simplify this to the form  $\left(\frac{x-A}{\alpha}\right)^2 + \left(\frac{y}{\beta}\right)^2 = 1$ . This is a nice exercise in messy algebra.

2. Let  $(x_1, y_1)$  and  $(x_2, y_2)$  be two points in  $\mathbb{R}^2$ . Give a simple description using the distance formula of the perpendicular bisector of the line segment joining these two points. Thus you want all points,  $(x, y)$  such that  $|(x, y) - (x_1, y_1)| = |(x, y) - (x_2, y_2)|$ .
3. Find a parametric equation for the line through the points  $(2, 3, 4, 5)$  and  $(-2, 3, 0, 1)$ .
4. Let  $(x, y) = (2 \cos(t), 2 \sin(t))$  where  $t \in [0, 2\pi]$ . Describe the set of points encountered as  $t$  changes.
5. Let  $(x, y, z) = (2 \cos(t), 2 \sin(t), t)$  where  $t \in \mathbb{R}$ . Describe the set of points encountered as  $t$  changes.

## 12.7 Open And Closed Sets

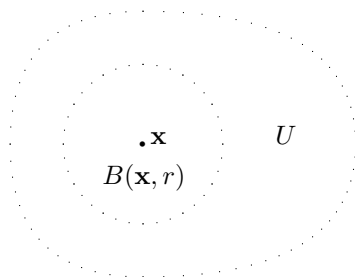
Eventually, one must consider functions which are defined on subsets of  $\mathbb{R}^n$  and their properties. The next definition will end up being quite important. It describes a type of subset of  $\mathbb{R}^n$  with the property that if  $\mathbf{x}$  is in this set, then so is  $\mathbf{y}$  whenever  $\mathbf{y}$  is close enough to  $\mathbf{x}$ .

**Definition 12.7.1** *Let  $U \subseteq \mathbb{R}^n$ .  $U$  is an open set if whenever  $\mathbf{x} \in U$ , there exists  $r > 0$  such that  $B(\mathbf{x}, r) \subseteq U$ . More generally, if  $U$  is any subset of  $\mathbb{R}^n$ ,  $\mathbf{x} \in U$  is an interior point of  $U$  if there exists  $r > 0$  such that  $\mathbf{x} \in B(\mathbf{x}, r) \subseteq U$ . In other words  $U$  is an open set exactly when every point of  $U$  is an interior point of  $U$ .*

If there is something called an open set, surely there should be something called a closed set and here is the definition of one.

**Definition 12.7.2** *A subset,  $C$ , of  $\mathbb{R}^n$  is called a closed set if  $\mathbb{R}^n \setminus C$  is an open set.*

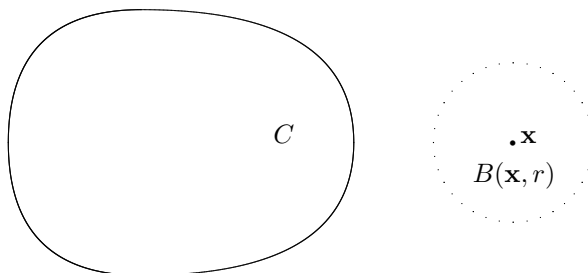
To illustrate this definition, consider the following picture.



You see in this picture how the edges are dotted. This is because an open set, can not include the edges or the set would fail to be open. For example, consider what would happen

if you picked a point out on the edge of  $U$  in the above picture. Every open ball centered at that point would have in it some points which are outside  $U$ . Therefore, such a point would violate the above definition. You also see the edges of  $B(\mathbf{x}, r)$  dotted suggesting that  $B(\mathbf{x}, r)$  ought to be an open set. This is intuitively clear but does require a proof. This will be done in the next theorem and will give examples of open sets. Also, you can see that if  $\mathbf{x}$  is close to the edge of  $U$ , you might have to take  $r$  to be very small.

It is roughly the case that open sets don't have their skins while closed sets do. Here is a picture of a closed set,  $C$ .



Note that  $\mathbf{x} \notin C$  and since  $\mathbb{R}^n \setminus C$  is open, there exists a ball,  $B(\mathbf{x}, r)$  contained entirely in  $\mathbb{R}^n \setminus C$ . If you look at  $\mathbb{R}^n \setminus C$ , what would be its skin? It can't be in  $\mathbb{R}^n \setminus C$  and so it must be in  $C$ . This is a rough heuristic explanation of what is going on with these definitions. Also note that  $\mathbb{R}^n$  and  $\emptyset$  are both open and closed. Here is why. If  $\mathbf{x} \in \emptyset$ , then there must be a ball centered at  $\mathbf{x}$  which is also contained in  $\emptyset$ . This must be considered to be true because there is nothing in  $\emptyset$  so there can be no example to show it false<sup>2</sup>. Therefore, from the definition, it follows  $\emptyset$  is open. It is also closed because if  $\mathbf{x} \notin \emptyset$ , then  $B(\mathbf{x}, 1)$  is also contained in  $\mathbb{R}^n \setminus \emptyset = \mathbb{R}^n$ . Therefore,  $\emptyset$  is both open and closed. From this, it follows  $\mathbb{R}^n$  is also both open and closed.

**Theorem 12.7.3** *Let  $\mathbf{x} \in \mathbb{R}^n$  and let  $r \geq 0$ . Then  $B(\mathbf{x}, r)$  is an open set. Also,*

$$D(\mathbf{x}, r) \equiv \{\mathbf{y} \in \mathbb{R}^n : |\mathbf{y} - \mathbf{x}| \leq r\}$$

*is a closed set.*

**Proof:** Suppose  $\mathbf{y} \in B(\mathbf{x}, r)$ . It is necessary to show there exists  $r_1 > 0$  such that  $B(\mathbf{y}, r_1) \subseteq B(\mathbf{x}, r)$ . Define  $r_1 \equiv r - |\mathbf{x} - \mathbf{y}|$ . Then if  $|\mathbf{z} - \mathbf{y}| < r_1$ , it follows from the above triangle inequality that

$$\begin{aligned} |\mathbf{z} - \mathbf{x}| &= |\mathbf{z} - \mathbf{y} + \mathbf{y} - \mathbf{x}| \\ &\leq |\mathbf{z} - \mathbf{y}| + |\mathbf{y} - \mathbf{x}| \\ &< r_1 + |\mathbf{y} - \mathbf{x}| = r - |\mathbf{x} - \mathbf{y}| + |\mathbf{y} - \mathbf{x}| = r. \end{aligned}$$

Note that if  $r = 0$  then  $B(\mathbf{x}, r) = \emptyset$ , the empty set. This is because if  $\mathbf{y} \in \mathbb{R}^n$ ,  $|\mathbf{x} - \mathbf{y}| \geq 0$  and so  $\mathbf{y} \notin B(\mathbf{x}, 0)$ . Since  $\emptyset$  has no points in it, it must be open because every point in it, (There are none.) satisfies the desired property of being an interior point.

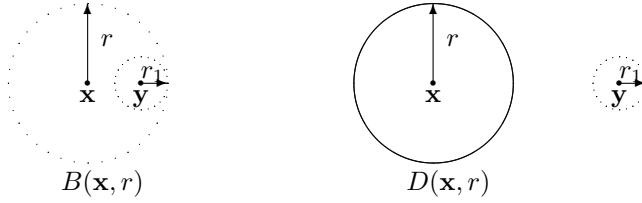
<sup>2</sup>To a mathematician, the statement: Whenever a pig is born with wings it can fly must be taken as true. We do not consider biological or aerodynamic considerations in such statements. There is no such thing as a winged pig and therefore, all winged pigs must be superb flyers since there can be no example of one which is not. On the other hand we would also consider the statement: Whenever a pig is born with wings it can't possibly fly, as equally true. The point is, you can say anything you want about the elements of the empty set and no one can gainsay your statement. Therefore, such statements are considered as true by default. You may say this is a very strange way of thinking about truth and ultimately this is because mathematics is not about truth. It is more about consistency and logic.

Now suppose  $\mathbf{y} \notin D(\mathbf{x}, r)$ . Then  $|\mathbf{x} - \mathbf{y}| > r$  and defining  $\delta \equiv |\mathbf{x} - \mathbf{y}| - r$ , it follows that if  $\mathbf{z} \in B(\mathbf{y}, \delta)$ , then by the triangle inequality,

$$\begin{aligned} |\mathbf{x} - \mathbf{z}| &\geq |\mathbf{x} - \mathbf{y}| - |\mathbf{y} - \mathbf{z}| > |\mathbf{x} - \mathbf{y}| - \delta \\ &= |\mathbf{x} - \mathbf{y}| - (|\mathbf{x} - \mathbf{y}| - r) = r \end{aligned}$$

and this shows that  $B(\mathbf{y}, \delta) \subseteq \mathbb{R}^n \setminus D(\mathbf{x}, r)$ . Since  $\mathbf{y}$  was an arbitrary point in  $\mathbb{R}^n \setminus D(\mathbf{x}, r)$ , it follows  $\mathbb{R}^n \setminus D(\mathbf{x}, r)$  is an open set which shows from the definition that  $D(\mathbf{x}, r)$  is a closed set as claimed.

A picture which is descriptive of the conclusion of the above theorem which also implies the manner of proof is the following.



Recall  $\mathbb{R}^2$  consists of ordered pairs,  $(x, y)$  such that  $x \in \mathbb{R}$  and  $y \in \mathbb{R}$ .  $\mathbb{R}^2$  is also written as  $\mathbb{R} \times \mathbb{R}$ . In general, the Cartesian product of two sets,  $A \times B$ , means  $\{(a, b) : a \in A, b \in B\}$ . Now suppose  $A \subseteq \mathbb{R}^m$  and  $B \subseteq \mathbb{R}^n$ . Then if  $(\mathbf{x}, \mathbf{y}) \in A \times B$ ,  $\mathbf{x} = (x_1, \dots, x_m)$  and  $\mathbf{y} = (y_1, \dots, y_n)$ . Therefore, the following identification will be made.

$$(\mathbf{x}, \mathbf{y}) = (x_1, \dots, x_m, y_1, \dots, y_n) \in \mathbb{R}^{n+m}.$$

Similarly, starting with something in  $\mathbb{R}^{n+m}$ , you can write it in the form  $(\mathbf{x}, \mathbf{y})$  where  $\mathbf{x} \in \mathbb{R}^m$  and  $\mathbf{y} \in \mathbb{R}^n$ . The following theorem has to do with the Cartesian product of two closed sets or two open sets. Also here is an important definition.

**Definition 12.7.4** A set,  $A \subseteq \mathbb{R}^n$  is said to be bounded if there exist finite intervals,  $[a_i, b_i]$  such that

$$A \subseteq \prod_{i=1}^n [a_i, b_i].$$

**Theorem 12.7.5** Let  $U$  be an open set in  $\mathbb{R}^m$  and let  $V$  be an open set in  $\mathbb{R}^n$ . Then  $U \times V$  is an open set in  $\mathbb{R}^{n+m}$ . If  $C$  is a closed set in  $\mathbb{R}^m$  and  $H$  is a closed set in  $\mathbb{R}^n$ , then  $C \times H$  is a closed set in  $\mathbb{R}^{n+m}$ . If  $C$  and  $H$  are bounded, then so is  $C \times H$ .

**Proof:** Let  $(\mathbf{x}, \mathbf{y}) \in U \times V$ . Since  $U$  is open, there exists  $r_1 > 0$  such that  $B(\mathbf{x}, r_1) \subseteq U$ . Similarly, there exists  $r_2 > 0$  such that  $B(\mathbf{y}, r_2) \subseteq V$ . Now

$$B((\mathbf{x}, \mathbf{y}), \delta) \equiv \left\{ (\mathbf{s}, \mathbf{t}) \in \mathbb{R}^{n+m} : \sum_{k=1}^m |x_k - s_k|^2 + \sum_{j=1}^n |y_j - t_j|^2 < \delta^2 \right\}$$

Therefore, if  $\delta \equiv \min(r_1, r_2)$  and  $(\mathbf{s}, \mathbf{t}) \in B((\mathbf{x}, \mathbf{y}), \delta)$ , then it follows that  $\mathbf{s} \in B(\mathbf{x}, r_1) \subseteq U$  and that  $\mathbf{t} \in B(\mathbf{y}, r_2) \subseteq V$  which shows that  $B((\mathbf{x}, \mathbf{y}), \delta) \subseteq U \times V$ . Hence  $U \times V$  is open as claimed.

Next suppose  $(\mathbf{x}, \mathbf{y}) \notin C \times H$ . It is necessary to show there exists  $\delta > 0$  such that  $B((\mathbf{x}, \mathbf{y}), \delta) \subseteq \mathbb{R}^{n+m} \setminus (C \times H)$ . Either  $\mathbf{x} \notin C$  or  $\mathbf{y} \notin H$  since otherwise  $(\mathbf{x}, \mathbf{y})$  would be a



point of  $C \times H$ . Suppose therefore, that  $\mathbf{x} \notin C$ . Since  $C$  is closed, there exists  $r > 0$  such that  $B(\mathbf{x}, r) \subseteq \mathbb{R}^m \setminus C$ . Consider  $B((\mathbf{x}, \mathbf{y}), r)$ . If  $(\mathbf{s}, \mathbf{t}) \in B((\mathbf{x}, \mathbf{y}), r)$ , it follows that  $\mathbf{s} \in B(\mathbf{x}, r)$  which is contained in  $\mathbb{R}^m \setminus C$ . Therefore,  $B((\mathbf{x}, \mathbf{y}), r) \subseteq \mathbb{R}^{n+m} \setminus (C \times H)$  showing  $C \times H$  is closed. A similar argument holds if  $\mathbf{y} \notin H$ .

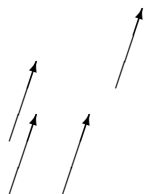
If  $C$  is bounded, there exist  $[a_i, b_i]$  such that  $C \subseteq \prod_{i=1}^m [a_i, b_i]$  and if  $H$  is bounded,  $H \subseteq \prod_{i=m+1}^{m+n} [a_i, b_i]$  for intervals  $[a_{m+1}, b_{m+1}], \dots, [a_{m+n}, b_{m+n}]$ . Therefore,  $C \times H \subseteq \prod_{i=1}^{m+n} [a_i, b_i]$  and this establishes the last part of this theorem.

## 12.8 Exercises

1. Show carefully that  $\mathbb{R}^n$  is both open and closed.
2. Show that every open set in  $\mathbb{R}^n$  is the union of open balls contained in it.
3. Show the intersection of any two open sets is an open set.

## 12.9 Vectors

Suppose you push on something. What is important? There are really two things which are important, how hard you push and the direction you push. Vectors are used to model this. What was just described would be called a force vector. It has two essential ingredients, its magnitude and its direction. Geometrically think of vectors as directed line segments as shown in the following picture in which all the directed line segments are considered to be the same vector because they have the same direction, the direction in which the arrows point, and the same magnitude (length).



Because of this fact that only direction and magnitude are important, it is always possible to put a vector in a certain particularly simple form. Let  $\overrightarrow{\mathbf{p}\mathbf{q}}$  be a directed line segment or vector. Then from Definition 12.5.3 that  $\overrightarrow{\mathbf{p}\mathbf{q}}$  consists of the points of the form

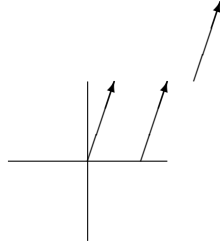
$$\mathbf{p} + t(\mathbf{q} - \mathbf{p})$$

where  $t \in [0, 1]$ . Subtract  $\mathbf{p}$  from all these points to obtain the directed line segment consisting of the points

$$\mathbf{0} + t(\mathbf{q} - \mathbf{p}), \quad t \in [0, 1].$$

The point in  $\mathbb{R}^n$ ,  $\mathbf{q} - \mathbf{p}$ , will represent the vector.

Geometrically, the arrow,  $\overrightarrow{\mathbf{p}\mathbf{q}}$ , was slid so it points in the same direction and the base is at the origin,  $\mathbf{0}$ . For example, see the following picture.



In this way vectors can be identified with elements of  $\mathbb{R}^n$ .

The magnitude of a vector determined by a directed line segment  $\overrightarrow{\mathbf{pq}}$  is just the distance between the point  $\mathbf{p}$  and the point  $\mathbf{q}$ . By the distance formula this equals

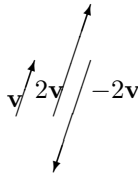
$$\left( \sum_{k=1}^n (q_k - p_k)^2 \right)^{1/2} = |\mathbf{p} - \mathbf{q}|$$

and for  $\mathbf{v}$  any vector in  $\mathbb{R}^n$  the magnitude of  $\mathbf{v}$  equals  $(\sum_{k=1}^n v_k^2)^{1/2} = |\mathbf{v}|$ .

What is the geometric significance of scalar multiplication? If  $\mathbf{a}$  represents the vector,  $\mathbf{v}$  in the sense that when it is slid to place its tail at the origin, the element of  $\mathbb{R}^n$  at its point is  $\mathbf{a}$ , what is  $r\mathbf{v}$ ?

$$\begin{aligned} |r\mathbf{v}| &= \left( \sum_{k=1}^n (ra_k)^2 \right)^{1/2} = \left( \sum_{k=1}^n r^2 (a_k)^2 \right)^{1/2} \\ &= (r^2)^{1/2} \left( \sum_{k=1}^n a_k^2 \right)^{1/2} = |r| |\mathbf{v}|. \end{aligned}$$

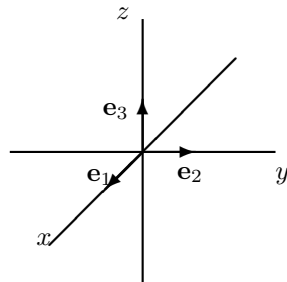
Thus the magnitude of  $r\mathbf{v}$  equals  $|r|$  times the magnitude of  $\mathbf{v}$ . If  $r$  is positive, then the vector represented by  $r\mathbf{v}$  has the same direction as the vector,  $\mathbf{v}$  because multiplying by the scalar,  $r$ , only has the effect of scaling all the distances. Thus the unit distance along any coordinate axis now has length  $r$  and in this rescaled system the vector is represented by  $\mathbf{a}$ . If  $r < 0$  similar considerations apply except in this case all the  $a_i$  also change sign. From now on,  $\mathbf{a}$  will be referred to as a vector instead of an element of  $\mathbb{R}^n$  representing a vector as just described. The following picture illustrates the effect of scalar multiplication.



Note there are  $n$  special vectors which point along the coordinate axes. These are

$$\mathbf{e}_i \equiv (0, \dots, 0, 1, 0, \dots, 0)$$

where the 1 is in the  $i^{th}$  slot and there are zeros in all the other spaces. See the picture in the case of  $\mathbb{R}^3$ .



The direction of  $\mathbf{e}_i$  is referred to as the  $i^{\text{th}}$  direction. Given a vector,  $\mathbf{v} = (a_1, \dots, a_n)$ ,  $a_i \mathbf{e}_i$  is the  $i^{\text{th}}$  component of the vector. Thus  $a_i \mathbf{e}_i = (0, \dots, 0, a_i, 0, \dots, 0)$  and so this vector gives something possibly nonzero only in the  $i^{\text{th}}$  direction. Also, knowledge of the  $i^{\text{th}}$  component of the vector is equivalent to knowledge of the vector because it gives the entry in the  $i^{\text{th}}$  slot and for  $\mathbf{v} = (a_1, \dots, a_n)$ ,

$$\mathbf{v} = \sum_{k=1}^n a_k \mathbf{e}_k.$$

What does addition of vectors mean physically? Suppose two forces are applied to some object. Each of these would be represented by a force vector and the two forces acting together would yield an overall force acting on the object which would also be a force vector known as the resultant. Suppose the two vectors are  $\mathbf{a} = \sum_{k=1}^n a_k \mathbf{e}_k$  and  $\mathbf{b} = \sum_{k=1}^n b_k \mathbf{e}_k$ . Then the vector,  $\mathbf{a}$  involves a component in the  $i^{\text{th}}$  direction,  $a_i \mathbf{e}_i$  while the component in the  $i^{\text{th}}$  direction of  $\mathbf{b}$  is  $b_i \mathbf{e}_i$ . Then it seems physically reasonable that the resultant vector should have a component in the  $i^{\text{th}}$  direction equal to  $(a_i + b_i) \mathbf{e}_i$ . This is exactly what is obtained when the vectors,  $\mathbf{a}$  and  $\mathbf{b}$  are added.

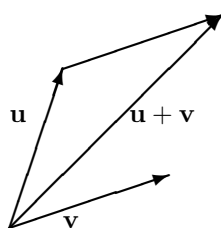
$$\begin{aligned} \mathbf{a} + \mathbf{b} &= (a_1 + b_1, \dots, a_n + b_n). \\ &= \sum_{i=1}^n (a_i + b_i) \mathbf{e}_i. \end{aligned}$$

Thus the addition of vectors according to the rules of addition in  $\mathbb{R}^n$ , yields the appropriate vector which duplicates the cumulative effect of all the vectors in the sum.

What is the geometric significance of vector addition? Suppose  $\mathbf{u}, \mathbf{v}$  are vectors,

$$\mathbf{u} = (u_1, \dots, u_n), \mathbf{v} = (v_1, \dots, v_n)$$

Then  $\mathbf{u} + \mathbf{v} = (u_1 + v_1, \dots, u_n + v_n)$ . How can one obtain this geometrically? Consider the directed line segment,  $\overrightarrow{0\mathbf{u}}$  and then, starting at the end of this directed line segment, follow the directed line segment  $\overrightarrow{\mathbf{u}(\mathbf{u} + \mathbf{v})}$  to its end,  $\mathbf{u} + \mathbf{v}$ . In other words, place the vector  $\mathbf{u}$  in standard position with its base at the origin and then slide the vector  $\mathbf{v}$  till its base coincides with the point of  $\mathbf{u}$ . The point of this slid vector, determines  $\mathbf{u} + \mathbf{v}$ . To illustrate, see the following picture



Note the vector  $\mathbf{u} + \mathbf{v}$  is the diagonal of a parallelogram determined from the two vectors  $\mathbf{u}$  and  $\mathbf{v}$  and that identifying  $\mathbf{u} + \mathbf{v}$  with the directed diagonal of the parallelogram determined by the vectors  $\mathbf{u}$  and  $\mathbf{v}$  amounts to the same thing as the above procedure.

An item of notation should be mentioned here. In the case of  $\mathbb{R}^n$  where  $n \leq 3$ , it is standard notation to use  $\mathbf{i}$  for  $\mathbf{e}_1$ ,  $\mathbf{j}$  for  $\mathbf{e}_2$ , and  $\mathbf{k}$  for  $\mathbf{e}_3$ . Now here are some applications of vector addition to some problems.

**Example 12.9.1** *There are three ropes attached to a car and three people pull on these ropes. The first exerts a force of  $2\mathbf{i} + 3\mathbf{j} - 2\mathbf{k}$  Newtons, the second exerts a force of  $3\mathbf{i} + 5\mathbf{j} + \mathbf{k}$*

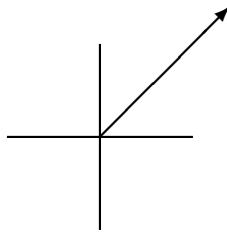
Newton and the third exerts a force of  $5\mathbf{i} - \mathbf{j} + 2\mathbf{k}$ . Newtons. Find the total force in the direction of  $\mathbf{i}$ .

To find the total force add the vectors as described above. This gives  $10\mathbf{i} + 7\mathbf{j} + \mathbf{k}$  Newtons. Therefore, the force in the  $\mathbf{i}$  direction is 10 Newtons.

The Newton is a unit of force like pounds.

**Example 12.9.2** An airplane flies North East at 100 miles per hour. Write this as a vector.

A picture of this situation follows.



The vector has length 100. Now using that vector as the hypotenuse of a right triangle having equal sides, the sides should be each of length  $100/\sqrt{2}$ . Therefore, the vector would be  $100/\sqrt{2}\mathbf{i} + 100/\sqrt{2}\mathbf{j}$ .

**Example 12.9.3** An airplane is traveling at  $100\mathbf{i} + \mathbf{j} + \mathbf{k}$  kilometers per hour and at a certain instant of time its position is  $(1, 2, 1)$ . Here imagine a Cartesian coordinate system in which the third component is altitude and the first and second components are measured on a line from West to East and a line from South to North. Find the position of this airplane one minute later.

Consider the vector  $(1, 2, 1)$ , is the initial position vector of the airplane. As it moves, the position vector changes. After one minute the airplane has moved in the  $\mathbf{i}$  direction a distance of  $100 \times \frac{1}{60} = \frac{5}{3}$  kilometer. In the  $\mathbf{j}$  direction it has moved  $\frac{1}{60}$  kilometer during this same time, while it moves  $\frac{1}{60}$  kilometer in the  $\mathbf{k}$  direction. Therefore, the new displacement vector for the airplane is

$$(1, 2, 1) + \left(\frac{5}{3}, \frac{1}{60}, \frac{1}{60}\right) = \left(\frac{8}{3}, \frac{121}{60}, \frac{121}{60}\right)$$

**Example 12.9.4** A certain river is one half mile wide with a current flowing at 4 miles per hour from East to West. A man swims directly toward the opposite shore from the South bank of the river at a speed of 3 miles per hour. How far down the river does he find himself when he has swam across? How far does he end up swimming?

Consider the following picture.



You should write these vectors in terms of components. The velocity of the swimmer in still water would be  $3\mathbf{j}$  while the velocity of the river would be  $-4\mathbf{i}$ . Therefore, the velocity

of the swimmer is  $-4\mathbf{i} + 3\mathbf{j}$ . Since the component of velocity in the direction across the river is 3, it follows the trip takes  $1/6$  hour or 10 minutes. The speed at which he travels is  $\sqrt{4^2 + 3^2} = 5$  miles per hour and so he travels  $5 \times \frac{1}{6} = \frac{5}{6}$  miles. Now to find the distance downstream he finds himself, note that if  $x$  is this distance,  $x$  and  $1/2$  are two legs of a right triangle whose hypotenuse equals  $5/6$  miles. Therefore, by the Pythagorean theorem the distance downstream is

$$\sqrt{(5/6)^2 - (1/2)^2} = \frac{2}{3} \text{ miles.}$$

## 12.10 Exercises

1. The wind blows from West to East at a speed of 50 kilometers per hour and an airplane is heading North West at a speed of 300 Kilometers per hour. What is the velocity of the airplane relative to the ground? What is the component of this velocity in the direction North?
2. In the situation of Problem 1 how many degrees to the West of North should the airplane head in order to fly exactly North. What will be the speed of the airplane?
3. In the situation of 2 suppose the airplane uses 34 gallons of fuel every hour at that air speed and that it needs to fly North a distance of 600 miles. Will the airplane have enough fuel to arrive at its destination given that it has 63 gallons of fuel?
4. A certain river is one half mile wide with a current flowing at 2 miles per hour from East to West. A man swims directly toward the opposite shore from the South bank of the river at a speed of 3 miles per hour. How far down the river does he find himself when he has swam across? How far does he end up swimming?
5. A certain river is one half mile wide with a current flowing at 2 miles per hour from East to West. A man can swim at 3 miles per hour in still water. In what direction should he swim in order to travel directly across the river? What would the answer to this problem be if the river flowed at 3 miles per hour and the man could swim only at the rate of 2 miles per hour?
6. Three forces are applied to a point which does not move. Two of the forces are  $2\mathbf{i} + \mathbf{j} + 3\mathbf{k}$  Newtons and  $\mathbf{i} - 3\mathbf{j} + 2\mathbf{k}$  Newtons. Find the third force.
7. An airplane is flying due north at 150 miles per hour. A wind is pushing the airplane due east at 40 miles per hour. After 1 hour, the plane starts flying  $30^\circ$  East of North. Assuming the plane starts at  $(0,0)$ , where is it after 2 hours? Let North be the direction of the positive  $y$  axis and let East be the direction of the positive  $x$  axis.
8. A bird flies from its nest 5 km. in the direction  $60^\circ$  north of east where it stops to rest on a tree. It then flies 10 km. in the direction due southeast and lands atop a telephone pole. Place an  $xy$  coordinate system so that the origin is the bird's nest, and the positive  $x$  axis points east and the positive  $y$  axis points north. Find the displacement vector from the nest to the telephone pole.



# Vector Products

## 13.1 The Dot Product

There are two ways of multiplying vectors which are of great importance in applications. The first of these is called the dot product, also called the scalar product and sometimes the inner product.

**Definition 13.1.1** Let  $\mathbf{a}, \mathbf{b}$  be two vectors in  $\mathbb{R}^n$  define  $\mathbf{a} \cdot \mathbf{b}$  as

$$\mathbf{a} \cdot \mathbf{b} \equiv \sum_{k=1}^n a_k b_k.$$

With this definition, there are several important properties satisfied by the dot product. In the statement of these properties,  $\alpha$  and  $\beta$  will denote scalars and  $\mathbf{a}, \mathbf{b}, \mathbf{c}$  will denote vectors.

**Proposition 13.1.2** The dot product satisfies the following properties.

$$\mathbf{a} \cdot \mathbf{b} = \mathbf{b} \cdot \mathbf{a} \quad (13.1)$$

$$\mathbf{a} \cdot \mathbf{a} \geq 0 \text{ and equals zero if and only if } \mathbf{a} = \mathbf{0} \quad (13.2)$$

$$(\alpha \mathbf{a} + \beta \mathbf{b}) \cdot \mathbf{c} = \alpha (\mathbf{a} \cdot \mathbf{c}) + \beta (\mathbf{b} \cdot \mathbf{c}) \quad (13.3)$$

$$\mathbf{c} \cdot (\alpha \mathbf{a} + \beta \mathbf{b}) = \alpha (\mathbf{c} \cdot \mathbf{a}) + \beta (\mathbf{c} \cdot \mathbf{b}) \quad (13.4)$$

$$|\mathbf{a}|^2 = \mathbf{a} \cdot \mathbf{a} \quad (13.5)$$

You should verify these properties. Also be sure you understand that (13.4) follows from the first three and is therefore redundant. It is listed here for the sake of convenience.

**Example 13.1.3** Find  $(1, 2, 0, -1) \cdot (0, 1, 2, 3)$ .

This equals  $0 + 2 + 0 + -3 = -1$ .

The dot product satisfies a fundamental inequality known as the Cauchy Schwartz inequality.

**Theorem 13.1.4** The dot product satisfies the inequality

$$|\mathbf{a} \cdot \mathbf{b}| \leq |\mathbf{a}| |\mathbf{b}|. \quad (13.6)$$

Furthermore equality is obtained if and only if one of  $\mathbf{a}$  or  $\mathbf{b}$  is a scalar multiple of the other.

**Proof:** First note that if  $\mathbf{b} = \mathbf{0}$  both sides of (13.6) equal zero and so the inequality holds in this case. Therefore, it will be assumed in what follows that  $\mathbf{b} \neq \mathbf{0}$ .

Define a function of  $t \in \mathbb{R}$

$$f(t) = (\mathbf{a} + t\mathbf{b}) \cdot (\mathbf{a} + t\mathbf{b}).$$

Then by (13.2),  $f(t) \geq 0$  for all  $t \in \mathbb{R}$ . Also from (13.3), (13.4), (13.1), and (13.5)

$$\begin{aligned} f(t) &= \mathbf{a} \cdot (\mathbf{a} + t\mathbf{b}) + t\mathbf{b} \cdot (\mathbf{a} + t\mathbf{b}) \\ &= \mathbf{a} \cdot \mathbf{a} + t(\mathbf{a} \cdot \mathbf{b}) + t\mathbf{b} \cdot \mathbf{a} + t^2\mathbf{b} \cdot \mathbf{b} \\ &= |\mathbf{a}|^2 + 2t(\mathbf{a} \cdot \mathbf{b}) + |\mathbf{b}|^2 t^2. \end{aligned}$$

Now

$$\begin{aligned} f(t) &= |\mathbf{b}|^2 \left( t^2 + 2t \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{b}|^2} + \frac{|\mathbf{a}|^2}{|\mathbf{b}|^2} \right) \\ &= |\mathbf{b}|^2 \left( t^2 + 2t \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{b}|^2} + \left( \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{b}|^2} \right)^2 - \left( \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{b}|^2} \right)^2 + \frac{|\mathbf{a}|^2}{|\mathbf{b}|^2} \right) \\ &= |\mathbf{b}|^2 \left( \left( t + \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{b}|^2} \right)^2 + \left( \frac{|\mathbf{a}|^2}{|\mathbf{b}|^2} - \left( \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{b}|^2} \right)^2 \right) \right) \geq 0 \end{aligned}$$

for all  $t \in \mathbb{R}$ . In particular  $f(t) \geq 0$  when  $t = -(\mathbf{a} \cdot \mathbf{b} / |\mathbf{b}|^2)$  which implies

$$\frac{|\mathbf{a}|^2}{|\mathbf{b}|^2} - \left( \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{b}|^2} \right)^2 \geq 0. \quad (13.7)$$

Multiplying both sides by  $|\mathbf{b}|^4$ ,

$$|\mathbf{a}|^2 |\mathbf{b}|^2 \geq (\mathbf{a} \cdot \mathbf{b})^2$$

which yields (13.6).

From Theorem 13.1.4, equality holds in (13.6) whenever one of the vectors is a scalar multiple of the other. It only remains to verify this is the only way equality can occur. If either vector equals zero, then equality is obtained in (13.6) so it can be assumed both vectors are non zero and that equality is obtained in (13.7). This implies that  $f(t) = 0$  when  $t = -(\mathbf{a} \cdot \mathbf{b} / |\mathbf{b}|^2)$  and so from (13.2), it follows that for this value of  $t$ ,  $\mathbf{a} + t\mathbf{b} = \mathbf{0}$  showing  $\mathbf{a} = -t\mathbf{b}$ . This proves the theorem.

You should note that the entire argument was based only on the properties of the dot product listed in (13.1) - (13.5). This means that whenever something satisfies these properties, the Cauchy Schwartz inequality holds. There are many other instances of these properties besides vectors in  $\mathbb{R}^n$ .

The Cauchy Schwartz inequality allows a proof of the triangle inequality for distances in  $\mathbb{R}^n$  in much the same way as the triangle inequality for the absolute value.

**Theorem 13.1.5** (*Triangle inequality*) For  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$

$$|\mathbf{a} + \mathbf{b}| \leq |\mathbf{a}| + |\mathbf{b}| \quad (13.8)$$

and equality holds if and only if one of the vectors is a nonnegative scalar multiple of the other. Also

$$||\mathbf{a}| - |\mathbf{b}|| \leq |\mathbf{a} - \mathbf{b}| \quad (13.9)$$



**Proof:** By properties of the dot product and the Cauchy Schwartz inequality,

$$\begin{aligned}
 |\mathbf{a} + \mathbf{b}|^2 &= (\mathbf{a} + \mathbf{b}) \cdot (\mathbf{a} + \mathbf{b}) \\
 &= (\mathbf{a} \cdot \mathbf{a}) + (\mathbf{a} \cdot \mathbf{b}) + (\mathbf{b} \cdot \mathbf{a}) + (\mathbf{b} \cdot \mathbf{b}) \\
 &= |\mathbf{a}|^2 + 2(\mathbf{a} \cdot \mathbf{b}) + |\mathbf{b}|^2 \\
 &\leq |\mathbf{a}|^2 + 2|\mathbf{a} \cdot \mathbf{b}| + |\mathbf{b}|^2 \\
 &\leq |\mathbf{a}|^2 + 2|\mathbf{a}||\mathbf{b}| + |\mathbf{b}|^2 \\
 &= (|\mathbf{a}| + |\mathbf{b}|)^2.
 \end{aligned}$$

Taking square roots of both sides you obtain (13.8).

It remains to consider when equality occurs. If either vector equals zero, then that vector equals zero times the other vector and the claim about when equality occurs is verified. Therefore, it can be assumed both vectors are nonzero. To get equality in the second inequality above, Theorem 13.1.4 implies one of the vectors must be a multiple of the other. Say  $\mathbf{b} = \alpha\mathbf{a}$ . If  $\alpha < 0$  then equality cannot occur in the first inequality because in this case

$$(\mathbf{a} \cdot \mathbf{b}) = \alpha |\mathbf{a}|^2 < 0 < |\alpha| |\mathbf{a}|^2 = |\mathbf{a} \cdot \mathbf{b}|$$

Therefore,  $\alpha \geq 0$ .

To get the other form of the triangle inequality,

$$\mathbf{a} = \mathbf{a} - \mathbf{b} + \mathbf{b}$$

so

$$\begin{aligned}
 |\mathbf{a}| &= |\mathbf{a} - \mathbf{b} + \mathbf{b}| \\
 &\leq |\mathbf{a} - \mathbf{b}| + |\mathbf{b}|.
 \end{aligned}$$

Therefore,

$$|\mathbf{a}| - |\mathbf{b}| \leq |\mathbf{a} - \mathbf{b}| \quad (13.10)$$

Similarly,

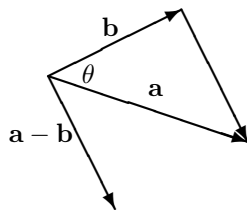
$$|\mathbf{b}| - |\mathbf{a}| \leq |\mathbf{b} - \mathbf{a}| = |\mathbf{a} - \mathbf{b}|. \quad (13.11)$$

It follows from (13.10) and (13.11) that (13.9) holds. This is because  $||\mathbf{a}| - |\mathbf{b}||$  equals the left side of either (13.10) or (13.11) and either way,  $||\mathbf{a}| - |\mathbf{b}|| \leq |\mathbf{a} - \mathbf{b}|$ . This proves the theorem.

## 13.2 The Geometric Significance Of The Dot Product

### 13.2.1 The Angle Between Two Vectors

Given two vectors,  $\mathbf{a}$  and  $\mathbf{b}$ , the included angle is the angle between these two vectors which is less than or equal to 180 degrees. The dot product can be used to determine the included angle between two vectors. To see how to do this, consider the following picture.



By the law of cosines,

$$|\mathbf{a} - \mathbf{b}|^2 = |\mathbf{a}|^2 + |\mathbf{b}|^2 - 2|\mathbf{a}||\mathbf{b}|\cos\theta.$$

Also from the properties of the dot product,

$$\begin{aligned} |\mathbf{a} - \mathbf{b}|^2 &= (\mathbf{a} - \mathbf{b}) \cdot (\mathbf{a} - \mathbf{b}) \\ &= |\mathbf{a}|^2 + |\mathbf{b}|^2 - 2\mathbf{a} \cdot \mathbf{b} \end{aligned}$$

and so comparing the above two formulas,

$$\mathbf{a} \cdot \mathbf{b} = |\mathbf{a}||\mathbf{b}|\cos\theta. \quad (13.12)$$

In words, the dot product of two vectors equals the product of the magnitude of the two vectors multiplied by the cosine of the included angle. Note this gives a geometric description of the dot product which does not depend explicitly on the coordinates of the vectors.

**Example 13.2.1** Find the angle between the vectors  $2\mathbf{i} + \mathbf{j} - \mathbf{k}$  and  $3\mathbf{i} + 4\mathbf{j} + \mathbf{k}$ .

The dot product of these two vectors equals  $6 + 4 - 1 = 9$  and the norms are  $\sqrt{4 + 1 + 1} = \sqrt{6}$  and  $\sqrt{9 + 16 + 1} = \sqrt{26}$ . Therefore, from (13.12) the cosine of the included angle equals

$$\cos\theta = \frac{9}{\sqrt{26}\sqrt{6}} = .72058$$

Now the cosine is known, the angle can be determined by solving the equation,  $\cos\theta = .72058$ . This will involve using a calculator or a table of trigonometric functions. The answer is  $\theta = .76616$  radians or in terms of degrees,  $\theta = .76616 \times \frac{360}{2\pi} = 43.898^\circ$ . Recall how this last computation is done. Set up a proportion,  $\frac{x}{.76616} = \frac{360}{2\pi}$  because  $360^\circ$  corresponds to  $2\pi$  radians. However, in calculus, you should get used to thinking in terms of radians and not degrees. This is because all the important calculus formulas are defined in terms of radians.

**Example 13.2.2** Let  $\mathbf{u}, \mathbf{v}$  be two vectors whose magnitudes are equal to 3 and 4 respectively and such that if they are placed in standard position with their tails at the origin, the angle between  $\mathbf{u}$  and the positive  $x$  axis equals  $30^\circ$  and the angle between  $\mathbf{v}$  and the positive  $x$  axis is  $-30^\circ$ . Find  $\mathbf{u} \cdot \mathbf{v}$ .

From the geometric description of the dot product in (13.12)

$$\mathbf{u} \cdot \mathbf{v} = 3 \times 4 \times \cos(60^\circ) = 3 \times 4 \times 1/2 = 6.$$

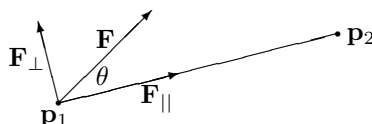
**Observation 13.2.3** Two vectors are said to be perpendicular if the included angle is  $\pi/2$  radians ( $90^\circ$ ). You can tell if two nonzero vectors are perpendicular by simply taking their dot product. If the answer is zero, this means they are perpendicular because  $\cos\theta = 0$ .

**Example 13.2.4** Determine whether the two vectors,  $2\mathbf{i} + \mathbf{j} - \mathbf{k}$  and  $\mathbf{i} + 3\mathbf{j} + 5\mathbf{k}$  are perpendicular.

When you take this dot product you get  $2 + 3 - 5 = 0$  and so these two are indeed perpendicular.

### 13.2.2 Work And Projections

Our first application will be to the concept of work. The physical concept of work does not in any way correspond to the notion of work employed in ordinary conversation. For example, if you were to slide a 150 pound weight off a table which is three feet high and shuffle along the floor for 50 yards, sweating profusely and exerting all your strength to keep the weight from falling on your feet, keeping the height always three feet and then deposit this weight on another three foot high table, the physical concept of work would indicate that the force exerted by your arms did no work during this project even though the muscles in your hands and arms would likely be very tired. The reason for such an unusual definition is that even though your arms exerted considerable force on the weight, enough to keep it from falling, the direction of motion was at right angles to the force they exerted. The only part of a force which does work in the sense of physics is the component of the force in the direction of motion. The work is defined to be the magnitude of the component of this force times the distance over which it acts in the case where this component of force points in the direction of motion and  $(-1)$  times the magnitude of this component times the distance in case the force tends to impede the motion. Thus the work done by a force on an object as the object moves from one point to another is a measure of the extent to which the force contributes to the motion. This is illustrated in the following picture in the case where the given force contributes to the motion.



In this picture the force,  $\mathbf{F}$  is applied to an object which moves on the straight line from  $\mathbf{p}_1$  to  $\mathbf{p}_2$ . There are two vectors shown,  $\mathbf{F}_{\parallel}$  and  $\mathbf{F}_{\perp}$  and the picture is intended to indicate that when you add these two vectors you get  $\mathbf{F}$  while  $\mathbf{F}_{\parallel}$  acts in the direction of motion and  $\mathbf{F}_{\perp}$  acts perpendicular to the direction of motion. Only  $\mathbf{F}_{\parallel}$  contributes to the work done by  $\mathbf{F}$  on the object as it moves from  $\mathbf{p}_1$  to  $\mathbf{p}_2$ . From trigonometry, you see the magnitude of  $\mathbf{F}_{\parallel}$  should equal  $|\mathbf{F}| |\cos \theta|$ . Thus, since  $\mathbf{F}_{\parallel}$  points in the direction of the vector from  $\mathbf{p}_1$  to  $\mathbf{p}_2$ , the total work done should equal

$$|\mathbf{F}| |\overrightarrow{\mathbf{p}_1 \mathbf{p}_2}| \cos \theta = |\mathbf{F}| |\mathbf{p}_2 - \mathbf{p}_1| \cos \theta$$

If the included angle had been obtuse, then the work done by the force,  $\mathbf{F}$  on the object would have been negative because in this case, the force tends to impede the motion from  $\mathbf{p}_1$  to  $\mathbf{p}_2$  but in this case,  $\cos \theta$  would also be negative and so it is still the case that the work done would be given by the above formula. Thus from the geometric description of the dot product given above, the work equals

$$|\mathbf{F}| |\mathbf{p}_2 - \mathbf{p}_1| \cos \theta = \mathbf{F} \cdot (\mathbf{p}_2 - \mathbf{p}_1).$$

This explains the following definition.

**Definition 13.2.5** *Let  $\mathbf{F}$  be a force acting on an object which moves from the point,  $\mathbf{p}_1$  to the point  $\mathbf{p}_2$ . Then the work done on the object by the given force equals  $\mathbf{F} \cdot (\mathbf{p}_2 - \mathbf{p}_1)$ .*

The concept of writing a given vector,  $\mathbf{F}$  in terms of two vectors, one which is parallel to a given vector,  $\mathbf{D}$  and the other which is perpendicular can also be explained with no reliance on trigonometry, completely in terms of the algebraic properties of the dot product.

As before, this is mathematically more significant than any approach involving geometry or trigonometry because it extends to more interesting situations. This is done next.

**Theorem 13.2.6** *Let  $\mathbf{F}$  and  $\mathbf{D}$  be nonzero vectors. Then there exist unique vectors  $\mathbf{F}_{||}$  and  $\mathbf{F}_{\perp}$  such that*

$$\mathbf{F} = \mathbf{F}_{||} + \mathbf{F}_{\perp} \quad (13.13)$$

where  $\mathbf{F}_{||}$  is a scalar multiple of  $\mathbf{D}$ , also referred to as

$$\text{proj}_{\mathbf{D}}(\mathbf{F}),$$

and  $\mathbf{F}_{\perp} \cdot \mathbf{D} = 0$ .

**Proof:** Suppose (13.13) and  $\mathbf{F}_{||} = \alpha \mathbf{D}$ . Taking the dot product of both sides with  $\mathbf{D}$  and using  $\mathbf{F}_{\perp} \cdot \mathbf{D} = 0$ , this yields

$$\mathbf{F} \cdot \mathbf{D} = \alpha |\mathbf{D}|^2$$

which requires  $\alpha = \mathbf{F} \cdot \mathbf{D} / |\mathbf{D}|^2$ . Thus there can be no more than one vector,  $\mathbf{F}_{||}$ . It follows  $\mathbf{F}_{\perp}$  must equal  $\mathbf{F} - \mathbf{F}_{||}$ . This verifies there can be no more than one choice for both  $\mathbf{F}_{||}$  and  $\mathbf{F}_{\perp}$ .

Now let

$$\mathbf{F}_{||} \equiv \frac{\mathbf{F} \cdot \mathbf{D}}{|\mathbf{D}|^2} \mathbf{D}$$

and let

$$\mathbf{F}_{\perp} = \mathbf{F} - \mathbf{F}_{||} = \mathbf{F} - \frac{\mathbf{F} \cdot \mathbf{D}}{|\mathbf{D}|^2} \mathbf{D}$$

Then  $\mathbf{F}_{||} = \alpha \mathbf{D}$  where  $\alpha = \frac{\mathbf{F} \cdot \mathbf{D}}{|\mathbf{D}|^2}$ . It only remains to verify  $\mathbf{F}_{\perp} \cdot \mathbf{D} = 0$ . But

$$\begin{aligned} \mathbf{F}_{\perp} \cdot \mathbf{D} &= \mathbf{F} \cdot \mathbf{D} - \frac{\mathbf{F} \cdot \mathbf{D}}{|\mathbf{D}|^2} \mathbf{D} \cdot \mathbf{D} \\ &= \mathbf{F} \cdot \mathbf{D} - \mathbf{F} \cdot \mathbf{D} = 0. \end{aligned}$$

This proves the theorem.

**Example 13.2.7** *Let  $\mathbf{F} = 2\mathbf{i} + 7\mathbf{j} - 3\mathbf{k}$  Newtons. Find the work done by this force in moving from the point  $(1, 2, 3)$  to the point  $(-9, -3, 4)$  where distances are measured in meters.*

According to the definition, this work is

$$\begin{aligned} (2\mathbf{i} + 7\mathbf{j} - 3\mathbf{k}) \cdot (-10\mathbf{i} - 5\mathbf{j} + \mathbf{k}) &= -20 + (-35) + (-3) \\ &= -58 \text{ Newton meters.} \end{aligned}$$

Note that if the force had been given in pounds and the distance had been given in feet, the units on the work would have been foot pounds. In general, work has units equal to units of a force times units of a length. Instead of writing Newton meter, people write joule because a joule is by definition a Newton meter. That word is pronounced “jewel” and it is the unit of work in the metric system of units. Also be sure you observe that the work done by the force can be negative as in the above example. In fact, work can be either positive, negative, or zero. You just have to do the computations to find out.

**Example 13.2.8** *Find  $\text{proj}_{\mathbf{u}}(\mathbf{v})$  if  $\mathbf{u} = 2\mathbf{i} + 3\mathbf{j} - 4\mathbf{k}$  and  $\mathbf{v} = \mathbf{i} - 2\mathbf{j} + \mathbf{k}$ .*

From the above discussion in Theorem 13.2.6, this is just

$$\begin{aligned} & \frac{1}{4+9+16} (\mathbf{i} - 2\mathbf{j} + \mathbf{k}) \cdot (2\mathbf{i} + 3\mathbf{j} - 4\mathbf{k}) (2\mathbf{i} + 3\mathbf{j} - 4\mathbf{k}) \\ &= \frac{-8}{29} (2\mathbf{i} + 3\mathbf{j} - 4\mathbf{k}) = -\frac{16}{29}\mathbf{i} - \frac{24}{29}\mathbf{j} + \frac{32}{29}\mathbf{k}. \end{aligned}$$

**Example 13.2.9** Suppose  $\mathbf{a}$ , and  $\mathbf{b}$  are vectors and  $\mathbf{b}_\perp = \mathbf{b} - \text{proj}_{\mathbf{a}}(\mathbf{b})$ . What is the magnitude of  $\mathbf{b}_\perp$  in terms of the included angle?

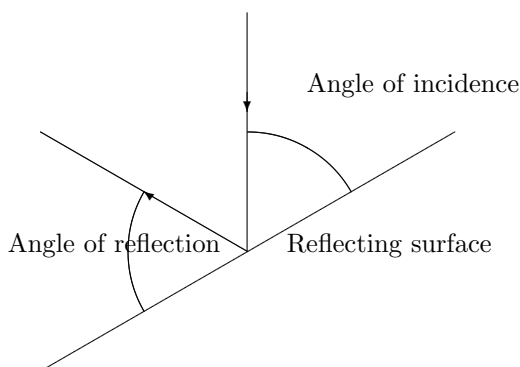
$$\begin{aligned} |\mathbf{b}_\perp|^2 &= (\mathbf{b} - \text{proj}_{\mathbf{a}}(\mathbf{b})) \cdot (\mathbf{b} - \text{proj}_{\mathbf{a}}(\mathbf{b})) \\ &= \left( \mathbf{b} - \frac{\mathbf{b} \cdot \mathbf{a}}{|\mathbf{a}|^2} \mathbf{a} \right) \cdot \left( \mathbf{b} - \frac{\mathbf{b} \cdot \mathbf{a}}{|\mathbf{a}|^2} \mathbf{a} \right) \\ &= |\mathbf{b}|^2 - 2 \frac{(\mathbf{b} \cdot \mathbf{a})^2}{|\mathbf{a}|^2} + \left( \frac{\mathbf{b} \cdot \mathbf{a}}{|\mathbf{a}|^2} \right)^2 |\mathbf{a}|^2 \\ &= |\mathbf{b}|^2 \left( 1 - \frac{(\mathbf{b} \cdot \mathbf{a})^2}{|\mathbf{a}|^2 |\mathbf{b}|^2} \right) \\ &= |\mathbf{b}|^2 (1 - \cos^2 \theta) = |\mathbf{b}|^2 \sin^2(\theta) \end{aligned}$$

where  $\theta$  is the included angle between  $\mathbf{a}$  and  $\mathbf{b}$  which is less than  $\pi$  radians. Therefore, taking square roots,

$$|\mathbf{b}_\perp| = |\mathbf{b}| \sin \theta.$$

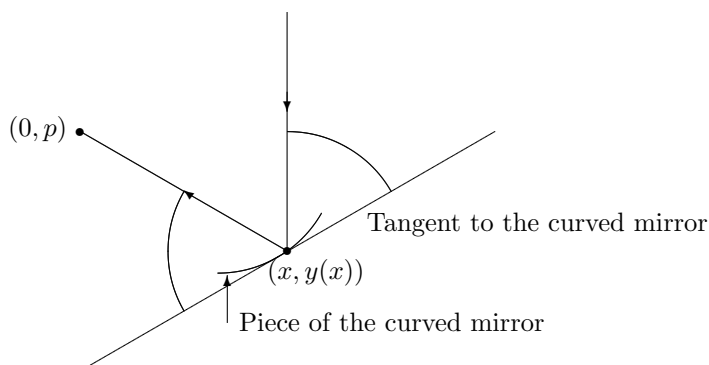
### 13.2.3 The Parabolic Mirror

When light is reflected the angle of incidence is always equal to the angle of reflection. This is illustrated in the following picture in which a ray of light reflects off something like a mirror.



An interesting problem is to design a curved mirror which has the property that it will direct all rays of light coming from a long distance away (essentially parallel rays of light) to a single point. You might be interested in a reflecting telescope for example or some sort of scheme for achieving high temperatures by reflecting the rays of the sun to a small area. Turning things around, you could place a source of light at the single point and desire to have the mirror reflect this in a beam of light consisting of parallel rays. How can you design such a mirror?

It turns out this is pretty easy given the above techniques for finding the angle between vectors. Consider the following picture.



It suffices to consider this in a plane for  $x > 0$  and then let the mirror be obtained as a surface of revolution. In the above picture, let  $(0, p)$  be the special point at which all the parallel rays of light will be directed. This is set up so the rays of light are parallel to the  $y$  axis. The two indicated angles will be equal and the equation of the indicated curve will be  $y = y(x)$  while the reflection is taking place at the point  $(x, y(x))$  as shown. To say the two angles are equal is to say their cosines are equal. Thus from the above,

$$\frac{(0, 1) \cdot (1, y'(x))}{\sqrt{1 + y'(x)^2}} = \frac{(-x, p - y) \cdot (-1, -y'(x))}{\sqrt{x^2 + (y - p)^2} \sqrt{1 + y'(x)^2}}.$$

This follows because the vectors forming the sides of one of the angles are  $(0, 1)$  and  $(1, y'(x))$  while the vectors forming the other angle are  $(-x, p - y)$  and  $(-1, -y'(x))$ . Therefore, this yields the differential equation,

$$y'(x) = \frac{-y'(x)(p - y) + x}{\sqrt{x^2 + (y - p)^2}}$$

which is written more simply as

$$\left( \sqrt{x^2 + (y - p)^2} + (p - y) \right) y' = x$$

Now let  $y - p = xv$  so that  $y' = xv' + v$ . Then in terms of  $v$  the differential equation is

$$xv' = \frac{1}{\sqrt{1 + v^2} - v} - v.$$

Using the technique for solving separable differential equations described in Problem 13 on Page 220 this reduces to

$$\left( \frac{1}{\sqrt{1 + v^2} - v} - v \right) dv = \frac{dx}{x}$$

To find  $\int \left( \frac{1}{\sqrt{1 + v^2} - v} - v \right) dv$  use a trig. substitution,  $v = \tan \theta$ . Then in terms of  $\theta$ , the antiderivative becomes

$$\begin{aligned} \int \left( \frac{1}{\sec \theta - \tan \theta} - \tan \theta \right) \sec^2 \theta d\theta &= \int \sec \theta d\theta \\ &= \ln |\sec \theta + \tan \theta| + C. \end{aligned}$$

Now in terms of  $v$ , this is

$$\ln(v + \sqrt{1 + v^2}) = \ln x + c.$$

There is no loss of generality in letting  $c = \ln C$  because  $\ln$  maps onto  $\mathbb{R}$ . Therefore, from laws of logarithms,

$$\begin{aligned} \ln|v + \sqrt{1 + v^2}| &= \ln x + c = \ln x + \ln C \\ &= \ln Cx. \end{aligned}$$

Therefore,

$$v + \sqrt{1 + v^2} = Cx$$

and so

$$\sqrt{1 + v^2} = Cx - v.$$

Now square both sides to get

$$1 + v^2 = C^2 x^2 + v^2 - 2C xv$$

which shows

$$1 = C^2 x^2 - 2Cx \frac{y - p}{x} = C^2 x^2 - 2C(y - p).$$

Solving this for  $y$  yields

$$y = \frac{C}{2}x^2 + \left(p - \frac{1}{2C}\right)$$

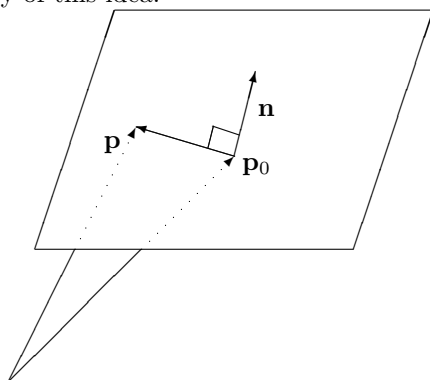
and for this to correspond to reflection as described above, it must be that  $C > 0$ . As described in an earlier section, this is just the equation of a parabola. Note it is possible to choose  $C$  as desired adjusting the shape of the mirror.

### 13.2.4 The Equation Of A Plane

The dot product makes possible a description of planes. To find the equation of a plane, you need two things, a point contained in the plane and a vector normal to the plane. Let  $\mathbf{p}_0 = (x_0, y_0, z_0)$  denote the position vector of a point in the plane, let  $\mathbf{p} = (x, y, z)$  be the position vector of an arbitrary point in the plane, and let  $\mathbf{n}$  denote a vector normal to the plane. This means that

$$\mathbf{n} \cdot (\mathbf{p} - \mathbf{p}_0) = 0$$

whenever  $\mathbf{p}$  is the position vector of a point in the plane. The following picture illustrates the geometry of this idea.



Expressed equivalently, the plane is just the set of all points  $\mathbf{p}$  such that the vector,  $\mathbf{p} - \mathbf{p}_0$  is perpendicular to the given normal vector,  $\mathbf{n}$ .

**Example 13.2.10** Find the equation of the plane with normal vector,  $\mathbf{n} = (1, 2, 3)$  containing the point  $(2, -1, 5)$ .

From the above, the equation of this plane is just

$$(1, 2, 3) \cdot (x - 2, y + 1, z - 3) = x - 9 + 2y + 3z = 0$$

**Example 13.2.11**  $2x + 4y - 5z = 11$  is the equation of a plane. Find the normal vector and a point on this plane.

You can write this in the form  $2(x - \frac{11}{2}) + 4(y - 0) + (-5)(z - 0) = 0$ . Therefore, a normal vector to the plane is  $2\mathbf{i} + 4\mathbf{j} - 5\mathbf{k}$  and a point in this plane is  $(\frac{11}{2}, 0, 0)$ . Of course there are many other points in the plane.

**Definition 13.2.12** Suppose two planes intersect. The angle between the planes is defined to be the angle between their normal vectors.

### 13.3 Exercises

1. Use formula (13.12) to verify the Cauchy Schwartz inequality and to show that equality occurs if and only if one of the vectors is a scalar multiple of the other.
2. For  $\mathbf{u}, \mathbf{v}$  vectors in  $\mathbb{R}^3$ , define the product,  $\mathbf{u} * \mathbf{v} \equiv u_1v_1 + 2u_2v_3 + 3u_3v_3$ . Show the axioms for a dot product all hold for this funny product. Prove  $|\mathbf{u} * \mathbf{v}| \leq (\mathbf{u} * \mathbf{u})^{1/2} (\mathbf{v} * \mathbf{v})^{1/2}$ .  
**Hint:** Do not try to do this with methods from trigonometry.
3. Find the angle between the vectors  $3\mathbf{i} - \mathbf{j} - \mathbf{k}$  and  $\mathbf{i} + 4\mathbf{j} + 2\mathbf{k}$ .
4. Find the angle between the vectors  $\mathbf{i} - 2\mathbf{j} + \mathbf{k}$  and  $\mathbf{i} + 2\mathbf{j} - 7\mathbf{k}$ .
5. If  $\mathbf{F}$  is a force and  $\mathbf{D}$  is a vector, show  $\text{proj}_{\mathbf{D}}(\mathbf{F}) = (|\mathbf{F}| \cos \theta) \mathbf{u}$  where  $\mathbf{u}$  is the unit vector in the direction of  $\mathbf{D}$ ,  $\mathbf{u} = \mathbf{D}/|\mathbf{D}|$  and  $\theta$  is the included angle between the two vectors,  $\mathbf{F}$  and  $\mathbf{D}$ .
6. A boy drags a sled for 100 feet along the ground by pulling on a rope which is 20 degrees from the horizontal with a force of 10 pounds. How much work does this force do?
7. An object moves 10 meters in the direction of  $\mathbf{j}$ . There are two forces acting on this object,  $\mathbf{F}_1 = \mathbf{i} + \mathbf{j} + 2\mathbf{k}$ , and  $\mathbf{F}_2 = -5\mathbf{i} + 2\mathbf{j} - 6\mathbf{k}$ . Find the total work done on the object by the two forces. **Hint:** You can take the work done by the resultant of the two forces or you can add the work done by each force.
8. An object moves 10 meters in the direction of  $\mathbf{j} + \mathbf{i}$ . There are two forces acting on this object,  $\mathbf{F}_1 = \mathbf{i} + 2\mathbf{j} + 2\mathbf{k}$ , and  $\mathbf{F}_2 = 5\mathbf{i} + 2\mathbf{j} - 6\mathbf{k}$ . Find the total work done on the object by the two forces. **Hint:** You can take the work done by the resultant of the two forces or you can add the work done by each force.
9. An object moves 20 meters in the direction of  $\mathbf{k} + \mathbf{j}$ . There are two forces acting on this object,  $\mathbf{F}_1 = \mathbf{i} + \mathbf{j} + 2\mathbf{k}$ , and  $\mathbf{F}_2 = \mathbf{i} + 2\mathbf{j} - 6\mathbf{k}$ . Find the total work done on the object by the two forces. **Hint:** You can take the work done by the resultant of the two forces or you can add the work done by each force.
10. If  $\mathbf{a}, \mathbf{b}$ , and  $\mathbf{c}$  are vectors. Show that  $(\mathbf{b} + \mathbf{c})_{\perp} = \mathbf{b}_{\perp} + \mathbf{c}_{\perp}$  where  $\mathbf{b}_{\perp} = \mathbf{b} - \text{proj}_{\mathbf{a}}(\mathbf{b})$ .



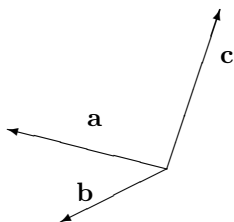
11. In the discussion of the reflecting mirror which directs all rays to a particular point,  $(0, p)$ . Show that for any choice of positive  $C$  this point is the focus of the parabola and the directrix is  $y = p - \frac{1}{C}$ .
12. Suppose you wanted to make a solar powered oven to cook food. Are there reasons for using a mirror which is not parabolic? Also describe how you would design a good flash light with a beam which does not spread out too quickly.
13. Find  $(1, 2, 3, 4) \cdot (2, 0, 1, 3)$ .
14. Show that  $(\mathbf{a} \cdot \mathbf{b}) = \frac{1}{4} \left[ |\mathbf{a} + \mathbf{b}|^2 - |\mathbf{a} - \mathbf{b}|^2 \right]$ .
15. Prove from the axioms of the dot product the parallelogram identity,  $|\mathbf{a} + \mathbf{b}|^2 + |\mathbf{a} - \mathbf{b}|^2 = 2|\mathbf{a}|^2 + 2|\mathbf{b}|^2$ .
16. Find the equation of the plane having a normal vector,  $5\mathbf{i} + 2\mathbf{j} - 6\mathbf{k}$  which contains the point  $(2, 1, 3)$ .
17. Find the equation of the plane having a normal vector,  $\mathbf{i} + 2\mathbf{j} - 4\mathbf{k}$  which contains the point  $(2, 0, 1)$ .
18. Find the equation of the plane having a normal vector,  $2\mathbf{i} + \mathbf{j} - 6\mathbf{k}$  which contains the point  $(1, 1, 2)$ .
19. Find the equation of the plane having a normal vector,  $\mathbf{i} + 2\mathbf{j} - 3\mathbf{k}$  which contains the point  $(1, 0, 3)$ .
20. If  $(a, b, c) \neq (0, 0, 0)$ , show that  $ax + by + cz = d$  is the equation of a plane with normal vector  $a\mathbf{i} + b\mathbf{j} + c\mathbf{k}$ .
21. Find the cosine of the angle between the two planes  $2x + 3y - z = 11$  and  $3x + y + 2z = 9$ .
22. Find the cosine of the angle between the two planes  $x + 3y - z = 11$  and  $2x + y + 2z = 9$ .
23. Find the cosine of the angle between the two planes  $2x + y - z = 11$  and  $3x + 5y + 2z = 9$ .
24. Find the cosine of the angle between the two planes  $x + 3y + z = 11$  and  $3x + 2y + 2z = 9$ .

## 13.4 The Cross Product

The cross product is the other way of multiplying two vectors in  $\mathbb{R}^3$ . It is very different from the dot product in many ways. First the geometric meaning is discussed and then a description in terms of coordinates is given. Both descriptions of the cross product are important. The geometric description is essential in order to understand the applications to physics and geometry while the coordinate description is the only way to practically compute the cross product.

**Definition 13.4.1** *Three vectors,  $\mathbf{a}, \mathbf{b}, \mathbf{c}$  form a right handed system if when you extend the fingers of your right hand along the vector,  $\mathbf{a}$  and close them in the direction of  $\mathbf{b}$ , the thumb points roughly in the direction of  $\mathbf{c}$ .*

For an example of a right handed system of vectors, see the following picture.



In this picture the vector  $\mathbf{c}$  points upwards from the plane determined by the other two vectors. You should consider how a right hand system would differ from a left hand system. Try using your left hand and you will see that the vector,  $\mathbf{c}$  would need to point in the opposite direction as it would for a right hand system.

From now on, the vectors,  $\mathbf{i}, \mathbf{j}, \mathbf{k}$  will always form a right handed system. To repeat, if you extend the fingers of our right hand along  $\mathbf{i}$  and close them in the direction  $\mathbf{j}$ , the thumb points in the direction of  $\mathbf{k}$ .

The following is the geometric description of the cross product. It gives both the direction and the magnitude and therefore specifies the vector.

**Definition 13.4.2** Let  $\mathbf{a}$  and  $\mathbf{b}$  be two vectors in  $\mathbb{R}^n$ . Then  $\mathbf{a} \times \mathbf{b}$  is defined by the following two rules.

1.  $|\mathbf{a} \times \mathbf{b}| = |\mathbf{a}| |\mathbf{b}| \sin \theta$  where  $\theta$  is the included angle.
2.  $\mathbf{a} \times \mathbf{b} \cdot \mathbf{a} = 0$ ,  $\mathbf{a} \times \mathbf{b} \cdot \mathbf{b} = 0$ , and  $\mathbf{a}, \mathbf{b}, \mathbf{a} \times \mathbf{b}$  forms a right hand system.

The cross product satisfies the following properties.

$$\mathbf{a} \times \mathbf{b} = -(\mathbf{b} \times \mathbf{a}), \quad \mathbf{a} \times \mathbf{a} = \mathbf{0}, \quad (13.14)$$

For  $\alpha$  a scalar,

$$(\alpha \mathbf{a}) \times \mathbf{b} = \alpha (\mathbf{a} \times \mathbf{b}) = \mathbf{a} \times (\alpha \mathbf{b}), \quad (13.15)$$

For  $\mathbf{a}, \mathbf{b}$ , and  $\mathbf{c}$  vectors, one obtains the distributive laws,

$$\mathbf{a} \times (\mathbf{b} + \mathbf{c}) = \mathbf{a} \times \mathbf{b} + \mathbf{a} \times \mathbf{c}, \quad (13.16)$$

$$(\mathbf{b} + \mathbf{c}) \times \mathbf{a} = \mathbf{b} \times \mathbf{a} + \mathbf{c} \times \mathbf{a}. \quad (13.17)$$

Formula (13.14) follows immediately from the definition. The vectors  $\mathbf{a} \times \mathbf{b}$  and  $\mathbf{b} \times \mathbf{a}$  have the same magnitude,  $|\mathbf{a}| |\mathbf{b}| \sin \theta$ , and an application of the right hand rule shows they have opposite direction. Formula (13.15) is also fairly clear. If  $\alpha$  is a nonnegative scalar, the direction of  $(\alpha \mathbf{a}) \times \mathbf{b}$  is the same as the direction of  $\mathbf{a} \times \mathbf{b}$ ,  $\alpha (\mathbf{a} \times \mathbf{b})$  and  $\mathbf{a} \times (\alpha \mathbf{b})$  while the magnitude is just  $\alpha$  times the magnitude of  $\mathbf{a} \times \mathbf{b}$  which is the same as the magnitude of  $\alpha (\mathbf{a} \times \mathbf{b})$  and  $\mathbf{a} \times (\alpha \mathbf{b})$ . Using this yields equality in (13.15). In the case where  $\alpha < 0$ , everything works the same way except the vectors are all pointing in the opposite direction and you must multiply by  $|\alpha|$  when comparing their magnitudes. The distributive laws are much harder to establish but the second follows from the first quite easily. Thus, assuming the first, and using (13.14),

$$\begin{aligned} (\mathbf{b} + \mathbf{c}) \times \mathbf{a} &= -\mathbf{a} \times (\mathbf{b} + \mathbf{c}) \\ &= -(\mathbf{a} \times \mathbf{b} + \mathbf{a} \times \mathbf{c}) \\ &= \mathbf{b} \times \mathbf{a} + \mathbf{c} \times \mathbf{a}. \end{aligned}$$

A proof of the distributive law is given in a later section for those who are interested. Now from the definition of the cross product,

$$\begin{aligned}\mathbf{i} \times \mathbf{j} &= \mathbf{k} & \mathbf{j} \times \mathbf{i} &= -\mathbf{k} \\ \mathbf{k} \times \mathbf{i} &= \mathbf{j} & \mathbf{i} \times \mathbf{k} &= -\mathbf{j} \\ \mathbf{j} \times \mathbf{k} &= \mathbf{i} & \mathbf{k} \times \mathbf{j} &= -\mathbf{i}\end{aligned}$$

With this information, the following gives the coordinate description of the cross product.

**Proposition 13.4.3** *Let  $\mathbf{a} = a_1\mathbf{i} + a_2\mathbf{j} + a_3\mathbf{k}$  and  $\mathbf{b} = b_1\mathbf{i} + b_2\mathbf{j} + b_3\mathbf{k}$  be two vectors. Then*

$$\begin{aligned}\mathbf{a} \times \mathbf{b} &= (a_2b_3 - a_3b_2)\mathbf{i} + (a_3b_1 - a_1b_3)\mathbf{j} + \\ &+ (a_1b_2 - a_2b_1)\mathbf{k}.\end{aligned}\tag{13.18}$$

**Proof:** From the above table and the properties of the cross product listed,

$$\begin{aligned} & (a_1\mathbf{i} + a_2\mathbf{j} + a_3\mathbf{k}) \times (b_1\mathbf{i} + b_2\mathbf{j} + b_3\mathbf{k}) = \\ & a_1b_2\mathbf{i} \times \mathbf{j} + a_1b_3\mathbf{i} \times \mathbf{k} + a_2b_1\mathbf{j} \times \mathbf{i} + a_2b_3\mathbf{j} \times \mathbf{k} + \\ & + a_3b_1\mathbf{k} \times \mathbf{i} + a_3b_2\mathbf{k} \times \mathbf{j} \\ & = a_1b_2\mathbf{k} - a_1b_3\mathbf{j} - a_2b_1\mathbf{k} + a_2b_3\mathbf{i} + a_3b_1\mathbf{j} - a_3b_2\mathbf{i} \\ & = (a_2b_3 - a_3b_2)\mathbf{i} + (a_3b_1 - a_1b_3)\mathbf{j} + (a_1b_2 - a_2b_1)\mathbf{k}\end{aligned}\tag{13.19}$$

This proves the proposition.

It is probably impossible for most people to remember (13.18). Fortunately, there is a somewhat easier way to remember it. This involves the notion of a determinant. A determinant is a single number assigned to a square array of numbers as follows.

$$\det \begin{pmatrix} a & b \\ c & d \end{pmatrix} \equiv \begin{vmatrix} a & b \\ c & d \end{vmatrix} = ad - bc.$$

This is the definition of the determinant of a square array of numbers having two rows and two columns. Now using this, the determinant of a square array of numbers in which there are three rows and three columns is defined as follows.

$$\begin{aligned}\det \begin{pmatrix} a & b & c \\ d & e & f \\ h & i & j \end{pmatrix} &\equiv (-1)^{1+1}a \begin{vmatrix} e & f \\ i & j \end{vmatrix} \\ &+ (-1)^{1+2}b \begin{vmatrix} d & f \\ h & j \end{vmatrix} + (-1)^{1+3}c \begin{vmatrix} d & e \\ h & i \end{vmatrix}.\end{aligned}$$

Take the first element in the top row,  $a$ , multiply by  $(-1)$  raised to the  $1 + 1$  since  $a$  is in the first row and the first column, and then multiply by the determinant obtained by crossing out the row and the column in which  $a$  appears. Then add to this a similar number obtained from the next element in the first row,  $b$ . This time multiply by  $(-1)^{1+2}$  because  $b$  is in the second column and the first row. When this is done do the same for  $c$ , the last element in the first row using a similar process. Using the definition of a determinant for square arrays of numbers having two columns and two rows, this equals

$$a(ej - if) + b(fh - dj) + c(di - eh),$$

an expression which, like the one for the cross product will be impossible to remember, although the process through which it is obtained is not too bad. It turns out these two

impossible to remember expressions are linked through the process of finding a determinant which was just described. The easy way to remember the description of the cross product in terms of coordinates, is to write

$$\mathbf{a} \times \mathbf{b} = \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{vmatrix} \quad (13.20)$$

and then follow the same process which was just described for calculating determinants above. This yields

$$(a_2b_3 - a_3b_2)\mathbf{i} - (a_1b_3 - a_3b_1)\mathbf{j} + (a_1b_2 - a_2b_1)\mathbf{k} \quad (13.21)$$

which is the same as (13.19). Later in the book a complete discussion of determinants is given but this will suffice for now.

**Example 13.4.4** Find  $(\mathbf{i} - \mathbf{j} + 2\mathbf{k}) \times (3\mathbf{i} - 2\mathbf{j} + \mathbf{k})$ .

Use (13.20) to compute this.

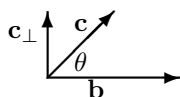
$$\begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ 1 & -1 & 2 \\ 3 & -2 & 1 \end{vmatrix} = \begin{vmatrix} -1 & 2 \\ -2 & 1 \end{vmatrix} \mathbf{i} - \begin{vmatrix} 1 & 2 \\ 3 & 1 \end{vmatrix} \mathbf{j} + \begin{vmatrix} 1 & -1 \\ 3 & -2 \end{vmatrix} \mathbf{k} \\ = 3\mathbf{i} + 5\mathbf{j} + \mathbf{k}.$$

### 13.4.1 The Distributive Law For The Cross Product

This section gives a proof for (13.16), a fairly difficult topic. It is included here for the interested student. If you are satisfied with taking the distributive law on faith, it is not necessary to read this section. The proof given here is quite clever and follows the one given in [5]. Another approach, based on volumes of parallelepipeds is found in [16] and is discussed a little later.

**Lemma 13.4.5** Let  $\mathbf{b}$  and  $\mathbf{c}$  be two vectors. Then  $\mathbf{b} \times \mathbf{c} = \mathbf{b} \times \mathbf{c}_\perp$  where  $\mathbf{c}_\parallel + \mathbf{c}_\perp = \mathbf{c}$  and  $\mathbf{c}_\perp \cdot \mathbf{b} = 0$ .

**Proof:** Consider the following picture.



Now  $\mathbf{c}_\perp = \mathbf{c} - \mathbf{c} \cdot \frac{\mathbf{b}}{|\mathbf{b}|} \frac{\mathbf{b}}{|\mathbf{b}|}$  and so  $\mathbf{c}_\perp$  is in the plane determined by  $\mathbf{c}$  and  $\mathbf{b}$ . Therefore, from the geometric definition of the cross product,  $\mathbf{b} \times \mathbf{c}$  and  $\mathbf{b} \times \mathbf{c}_\perp$  have the same direction. Now, referring to the picture,

$$\begin{aligned} |\mathbf{b} \times \mathbf{c}_\perp| &= |\mathbf{b}| |\mathbf{c}_\perp| \\ &= |\mathbf{b}| |\mathbf{c}| \sin \theta \\ &= |\mathbf{b} \times \mathbf{c}|. \end{aligned}$$

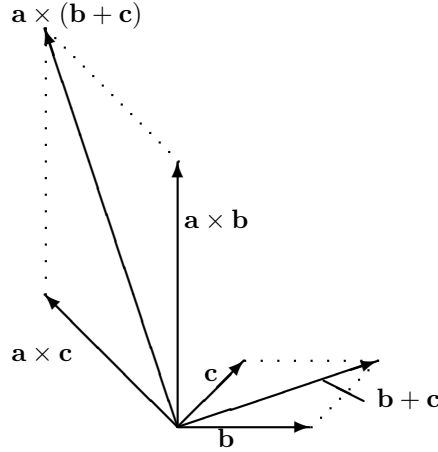
Therefore,  $\mathbf{b} \times \mathbf{c}$  and  $\mathbf{b} \times \mathbf{c}_\perp$  also have the same magnitude and so they are the same vector.

With this, the proof of the distributive law is in the following theorem.

**Theorem 13.4.6** Let  $\mathbf{a}, \mathbf{b}$ , and  $\mathbf{c}$  be vectors in  $\mathbb{R}^3$ . Then

$$\mathbf{a} \times (\mathbf{b} + \mathbf{c}) = \mathbf{a} \times \mathbf{b} + \mathbf{a} \times \mathbf{c} \quad (13.22)$$

**Proof:** Suppose first that  $\mathbf{a} \cdot \mathbf{b} = \mathbf{a} \cdot \mathbf{c} = 0$ . Now imagine  $\mathbf{a}$  is a vector coming out of the page and let  $\mathbf{b}, \mathbf{c}$  and  $\mathbf{b} + \mathbf{c}$  be as shown in the following picture.



Then  $\mathbf{a} \times \mathbf{b}, \mathbf{a} \times (\mathbf{b} + \mathbf{c})$ , and  $\mathbf{a} \times \mathbf{c}$  are each vectors in the same plane, perpendicular to  $\mathbf{a}$  as shown. Thus  $\mathbf{a} \times \mathbf{c} \cdot \mathbf{c} = 0, \mathbf{a} \times (\mathbf{b} + \mathbf{c}) \cdot (\mathbf{b} + \mathbf{c}) = 0$ , and  $\mathbf{a} \times \mathbf{b} \cdot \mathbf{b} = 0$ . This implies that to get  $\mathbf{a} \times \mathbf{b}$  you move counterclockwise through an angle of  $\pi/2$  radians from the vector,  $\mathbf{b}$ . Similar relationships exist between the vectors  $\mathbf{a} \times (\mathbf{b} + \mathbf{c})$  and  $\mathbf{b} + \mathbf{c}$  and the vectors  $\mathbf{a} \times \mathbf{c}$  and  $\mathbf{c}$ . Thus the angle between  $\mathbf{a} \times \mathbf{b}$  and  $\mathbf{a} \times (\mathbf{b} + \mathbf{c})$  is the same as the angle between  $\mathbf{b} + \mathbf{c}$  and  $\mathbf{b}$  and the angle between  $\mathbf{a} \times \mathbf{c}$  and  $\mathbf{a} \times (\mathbf{b} + \mathbf{c})$  is the same as the angle between  $\mathbf{c}$  and  $\mathbf{b} + \mathbf{c}$ . In addition to this, since  $\mathbf{a}$  is perpendicular to these vectors,

$$|\mathbf{a} \times \mathbf{b}| = |\mathbf{a}| |\mathbf{b}|, |\mathbf{a} \times (\mathbf{b} + \mathbf{c})| = |\mathbf{a}| |\mathbf{b} + \mathbf{c}|, \text{ and}$$

$$|\mathbf{a} \times \mathbf{c}| = |\mathbf{a}| |\mathbf{c}|.$$

Therefore,

$$\frac{|\mathbf{a} \times (\mathbf{b} + \mathbf{c})|}{|\mathbf{b} + \mathbf{c}|} = \frac{|\mathbf{a} \times \mathbf{c}|}{|\mathbf{c}|} = \frac{|\mathbf{a} \times \mathbf{b}|}{|\mathbf{b}|} = |\mathbf{a}|$$

and so

$$\frac{|\mathbf{a} \times (\mathbf{b} + \mathbf{c})|}{|\mathbf{a} \times \mathbf{c}|} = \frac{|\mathbf{b} + \mathbf{c}|}{|\mathbf{c}|}, \quad \frac{|\mathbf{a} \times (\mathbf{b} + \mathbf{c})|}{|\mathbf{a} \times \mathbf{b}|} = \frac{|\mathbf{b} + \mathbf{c}|}{|\mathbf{b}|}$$

showing the triangles making up the parallelogram on the right and the four sided figure on the left in the above picture are similar. It follows the four sided figure on the left is in fact a parallelogram and this implies the diagonal is the vector sum of the vectors on the sides, yielding (13.22).

Now suppose it is not necessarily the case that  $\mathbf{a} \cdot \mathbf{b} = \mathbf{a} \cdot \mathbf{c} = 0$ . Then write  $\mathbf{b} = \mathbf{b}_{\parallel} + \mathbf{b}_{\perp}$  where  $\mathbf{b}_{\perp} \cdot \mathbf{a} = 0$ . Similarly  $\mathbf{c} = \mathbf{c}_{\parallel} + \mathbf{c}_{\perp}$ . By the above lemma and what was just shown,

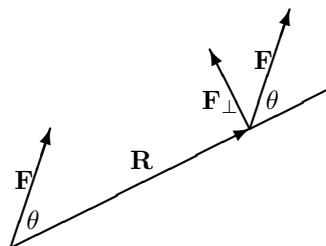
$$\begin{aligned} \mathbf{a} \times (\mathbf{b} + \mathbf{c}) &= \mathbf{a} \times (\mathbf{b} + \mathbf{c})_{\perp} \\ &= \mathbf{a} \times (\mathbf{b}_{\perp} + \mathbf{c}_{\perp}) \\ &= \mathbf{a} \times \mathbf{b}_{\perp} + \mathbf{a} \times \mathbf{c}_{\perp} \\ &= \mathbf{a} \times \mathbf{b} + \mathbf{a} \times \mathbf{c}. \end{aligned}$$

This proves the theorem.

The result of Problem 10 of the exercises 13.3 is used to go from the first to the second line.

### 13.4.2 Torque

Imagine you are using a wrench to loosen a nut. The idea is to turn the nut by applying a force to the end of the wrench. If you push or pull the wrench directly toward or away from the nut, it should be obvious from experience that no progress will be made in turning the nut. The important thing is the component of force perpendicular to the wrench. It is this component of force which will cause the nut to turn. For example see the following picture.



In the picture a force,  $\mathbf{F}$  is applied at the end of a wrench represented by the position vector,  $\mathbf{R}$  and the angle between these two is  $\theta$ . Then the tendency to turn will be  $|\mathbf{R}| |\mathbf{F}_\perp| = |\mathbf{R}| |\mathbf{F}| \sin \theta$ , which you recognize as the magnitude of the cross product of  $\mathbf{R}$  and  $\mathbf{F}$ . If there were just one force acting at one point whose position vector is  $\mathbf{R}$ , perhaps this would be sufficient, but what if there are numerous forces acting at many different points with neither the position vectors nor the force vectors in the same plane; what then? To keep track of this sort of thing, define for each  $\mathbf{R}$  and  $\mathbf{F}$ , the Torque vector,

$$\tau \equiv \mathbf{R} \times \mathbf{F}.$$

That way, if there are several forces acting at several points, the total torque can be obtained by simply adding up the torques associated with the different forces and positions.

**Example 13.4.7** Suppose  $\mathbf{R}_1 = 2\mathbf{i} - \mathbf{j} + 3\mathbf{k}$ ,  $\mathbf{R}_2 = \mathbf{i} + 2\mathbf{j} - 6\mathbf{k}$  meters and at the points determined by these vectors there are forces,  $\mathbf{F}_1 = \mathbf{i} - \mathbf{j} + 2\mathbf{k}$  and  $\mathbf{F}_2 = \mathbf{i} - 5\mathbf{j} + \mathbf{k}$  Newtons respectively. Find the total torque about the origin produced by these forces acting at the given points.

It is necessary to take  $\mathbf{R}_1 \times \mathbf{F}_1 + \mathbf{R}_2 \times \mathbf{F}_2$ . Thus the total torque equals

$$\begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ 2 & -1 & 3 \\ 1 & -1 & 2 \end{vmatrix} + \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ 1 & 2 & -6 \\ 1 & -5 & 1 \end{vmatrix} = -27\mathbf{i} - 8\mathbf{j} - 8\mathbf{k} \text{ Newton meters}$$

**Example 13.4.8** Find if possible a single force vector,  $\mathbf{F}$  which if applied at the point  $\mathbf{i} + \mathbf{j} + \mathbf{k}$  will produce the same torque as the above two forces acting at the given points.

This is fairly routine. The problem is to find  $\mathbf{F} = F_1\mathbf{i} + F_2\mathbf{j} + F_3\mathbf{k}$  which produces the above torque vector. Therefore,

$$\begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ 1 & 1 & 1 \\ F_1 & F_2 & F_3 \end{vmatrix} = -27\mathbf{i} - 8\mathbf{j} - 8\mathbf{k}$$

which reduces to  $(F_3 - F_2)\mathbf{i} + (F_1 - F_3)\mathbf{j} + (F_2 - F_1)\mathbf{k} = -27\mathbf{i} - 8\mathbf{j} - 8\mathbf{k}$ . This amounts to solving the system of three equations in three unknowns,  $F_1, F_2$ , and  $F_3$ ,

$$\begin{aligned}F_3 - F_2 &= -27 \\F_1 - F_3 &= -8 \\F_2 - F_1 &= -8\end{aligned}$$

However, there is no solution to these three equations. (Why?) Therefore no single force acting at the point  $\mathbf{i} + \mathbf{j} + \mathbf{k}$  will produce the given torque.

The mass of an object is a measure of how much stuff there is in the object. An object has mass equal to one kilogram, a unit of mass in the metric system, if it would exactly balance a known one kilogram object when placed on a balance. The known object is one kilogram by definition. The mass of an object does not depend on where the balance is used. It would be one kilogram on the moon as well as on the earth. The weight of an object is something else. It is the force exerted on the object by gravity and has magnitude  $gm$  where  $g$  is a constant called the acceleration of gravity. Thus the weight of a one kilogram object would be different on the moon which has much less gravity, smaller  $g$ , than on the earth. An important idea is that of the center of mass. This is the point at which an object will balance no matter how it is turned.

**Definition 13.4.9** Let an object consist of  $p$  point masses,  $m_1, \dots, m_p$  with the position of the  $k^{\text{th}}$  of these at  $\mathbf{R}_k$ . The center of mass of this object,  $\mathbf{R}_0$  is the point satisfying

$$\sum_{k=1}^p (\mathbf{R}_k - \mathbf{R}_0) \times gm_k \mathbf{u} = \mathbf{0}$$

for all unit vectors,  $\mathbf{u}$ .

The above definition indicates that no matter how the object is suspended, the total torque on it due to gravity is such that no rotation occurs. Using the properties of the cross product,

$$\left( \sum_{k=1}^p \mathbf{R}_k gm_k - \mathbf{R}_0 \sum_{k=1}^p gm_k \right) \times \mathbf{u} = \mathbf{0} \quad (13.23)$$

for any choice of unit vector,  $\mathbf{u}$ . You should verify that if  $\mathbf{a} \times \mathbf{u} = \mathbf{0}$  for all  $\mathbf{u}$ , then it must be the case that  $\mathbf{a} = \mathbf{0}$ . Then the above formula requires that

$$\sum_{k=1}^p \mathbf{R}_k gm_k - \mathbf{R}_0 \sum_{k=1}^p gm_k = \mathbf{0}.$$

dividing by  $g$ , and then by  $\sum_{k=1}^p m_k$ ,

$$\mathbf{R}_0 = \frac{\sum_{k=1}^p \mathbf{R}_k m_k}{\sum_{k=1}^p m_k}. \quad (13.24)$$

This is the formula for the center of mass of a collection of point masses. To consider the center of mass of a solid consisting of continuously distributed masses, you need the methods of calculus.

**Example 13.4.10** Let  $m_1 = 5, m_2 = 6$ , and  $m_3 = 3$  where the masses are in kilograms. Suppose  $m_1$  is located at  $2\mathbf{i} + 3\mathbf{j} + \mathbf{k}$ ,  $m_2$  is located at  $\mathbf{i} - 3\mathbf{j} + 2\mathbf{k}$  and  $m_3$  is located at  $2\mathbf{i} - \mathbf{j} + 3\mathbf{k}$ . Find the center of mass of these three masses.

Using (13.24)

$$\begin{aligned}\mathbf{R}_0 &= \frac{5(2\mathbf{i} + 3\mathbf{j} + \mathbf{k}) + 6(\mathbf{i} - 3\mathbf{j} + 2\mathbf{k}) + 3(2\mathbf{i} - \mathbf{j} + 3\mathbf{k})}{5 + 6 + 3} \\ &= \frac{11}{7}\mathbf{i} - \frac{3}{7}\mathbf{j} + \frac{13}{7}\mathbf{k}\end{aligned}$$

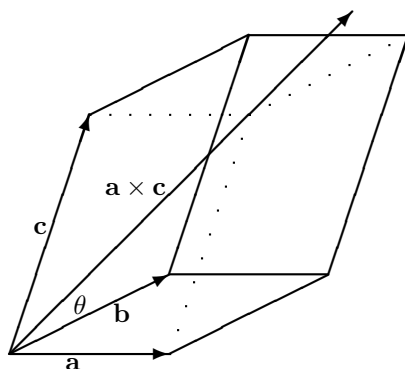
### 13.4.3 The Box Product

**Definition 13.4.11** A parallelepiped determined by the three vectors,  $\mathbf{a}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$  consists of

$$\{r\mathbf{a} + s\mathbf{b} + t\mathbf{c} : r, s, t \in [0, 1]\}.$$

That is, if you pick three numbers,  $r$ ,  $s$ , and  $t$  each in  $[0, 1]$  and form  $r\mathbf{a} + s\mathbf{b} + t\mathbf{c}$ , then the collection of all such points is what is meant by the parallelepiped determined by these three vectors.

The following is a picture of such a thing.



You notice the area of the base of the parallelepiped, the parallelogram determined by the vectors,  $\mathbf{a}$  and  $\mathbf{c}$  has area equal to  $|\mathbf{a} \times \mathbf{c}|$  while the altitude of the parallelepiped is  $|\mathbf{b}| \cos \theta$  where  $\theta$  is the angle shown in the picture between  $\mathbf{b}$  and  $\mathbf{a} \times \mathbf{c}$ . Therefore, the volume of this parallelepiped is the area of the base times the altitude which is just

$$|\mathbf{a} \times \mathbf{c}| |\mathbf{b}| \cos \theta = \mathbf{a} \times \mathbf{c} \cdot \mathbf{b}.$$

This expression is known as the box product and is sometimes written as  $[\mathbf{a}, \mathbf{c}, \mathbf{b}]$ . You should consider what happens if you interchange the  $\mathbf{b}$  with the  $\mathbf{c}$  or the  $\mathbf{a}$  with the  $\mathbf{c}$ . You can see geometrically from drawing pictures that this merely introduces a minus sign. In any case the box product of three vectors always equals either the volume of the parallelepiped determined by the three vectors or else minus this volume.

**Example 13.4.12** Find the volume of the parallelepiped determined by the vectors,  $\mathbf{i} + 2\mathbf{j} - 5\mathbf{k}$ ,  $\mathbf{i} + 3\mathbf{j} - 6\mathbf{k}$ ,  $3\mathbf{i} + 2\mathbf{j} + 3\mathbf{k}$ .

According to the above discussion, pick any two of these, take the cross product and then take the dot product of this with the third of these vectors. The result will be either the desired volume or minus the desired volume.

$$\begin{aligned}(\mathbf{i} + 2\mathbf{j} - 5\mathbf{k}) \times (\mathbf{i} + 3\mathbf{j} - 6\mathbf{k}) &= \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ 1 & 2 & -5 \\ 1 & 3 & -6 \end{vmatrix} \\ &= 3\mathbf{i} + \mathbf{j} + \mathbf{k}\end{aligned}$$



Now take the dot product of this vector with the third which yields

$$(3\mathbf{i} + \mathbf{j} + \mathbf{k}) \cdot (3\mathbf{i} + 2\mathbf{j} + 3\mathbf{k}) = 9 + 2 + 3 = 14.$$

This shows the volume of this parallelepiped is 14 cubic units.

Here is another proof of the distributive law for the cross product. From the above picture  $(\mathbf{a} \times \mathbf{b}) \cdot \mathbf{c} = \mathbf{a} \cdot (\mathbf{b} \times \mathbf{c})$  because both of these give either the volume of a parallelepiped determined by the vectors  $\mathbf{a}, \mathbf{b}, \mathbf{c}$  or -1 times the volume of this parallelepiped. Now to prove the distributive law, let  $\mathbf{x}$  be a vector. From the above observation,

$$\begin{aligned} \mathbf{x} \cdot \mathbf{a} \times (\mathbf{b} + \mathbf{c}) &= (\mathbf{x} \times \mathbf{a}) \cdot (\mathbf{b} + \mathbf{c}) \\ &= (\mathbf{x} \times \mathbf{a}) \cdot \mathbf{b} + (\mathbf{x} \times \mathbf{a}) \cdot \mathbf{c} \\ &= \mathbf{x} \cdot \mathbf{a} \times \mathbf{b} + \mathbf{x} \cdot \mathbf{a} \times \mathbf{c} \\ &= \mathbf{x} \cdot (\mathbf{a} \times \mathbf{b} + \mathbf{a} \times \mathbf{c}). \end{aligned}$$

Therefore,

$$\mathbf{x} \cdot [\mathbf{a} \times (\mathbf{b} + \mathbf{c}) - (\mathbf{a} \times \mathbf{b} + \mathbf{a} \times \mathbf{c})] = 0$$

for all  $\mathbf{x}$ . In particular, this holds for  $\mathbf{x} = \mathbf{a} \times (\mathbf{b} + \mathbf{c}) - (\mathbf{a} \times \mathbf{b} + \mathbf{a} \times \mathbf{c})$  showing that  $\mathbf{a} \times (\mathbf{b} + \mathbf{c}) = \mathbf{a} \times \mathbf{b} + \mathbf{a} \times \mathbf{c}$  and this proves the distributive law for the cross product another way.

## 13.5 Exercises

1. Show that if  $\mathbf{a} \times \mathbf{u} = \mathbf{0}$  for all unit vectors,  $\mathbf{u}$ , then  $\mathbf{a} = \mathbf{0}$ .
2. If you only assume (13.23) holds for  $\mathbf{u} = \mathbf{i}, \mathbf{j}, \mathbf{k}$ , show that this implies (13.23) holds for all unit vectors,  $\mathbf{u}$ .
3. Let  $m_1 = 5, m_2 = 1$ , and  $m_3 = 4$  where the masses are in kilograms and the distance is in meters. Suppose  $m_1$  is located at  $2\mathbf{i} - 3\mathbf{j} + \mathbf{k}$ ,  $m_2$  is located at  $\mathbf{i} - 3\mathbf{j} + 6\mathbf{k}$  and  $m_3$  is located at  $2\mathbf{i} + \mathbf{j} + 3\mathbf{k}$ . Find the center of mass of these three masses.
4. Let  $m_1 = 2, m_2 = 3$ , and  $m_3 = 1$  where the masses are in kilograms and the distance is in meters. Suppose  $m_1$  is located at  $2\mathbf{i} - \mathbf{j} + \mathbf{k}$ ,  $m_2$  is located at  $\mathbf{i} - 2\mathbf{j} + \mathbf{k}$  and  $m_3$  is located at  $4\mathbf{i} + \mathbf{j} + 3\mathbf{k}$ . Find the center of mass of these three masses.
5. Find the volume of the parallelepiped determined by the vectors,  $\mathbf{i} - 7\mathbf{j} - 5\mathbf{k}, \mathbf{i} - 2\mathbf{j} - 6\mathbf{k}, 3\mathbf{i} + 2\mathbf{j} + 3\mathbf{k}$ .
6. Suppose  $\mathbf{a}, \mathbf{b}$ , and  $\mathbf{c}$  are three vectors whose components are all integers. Can you conclude the volume of the parallelepiped determined from these three vectors will always be an integer?
7. What does it mean geometrically if the box product of three vectors gives zero?
8. Suppose  $\mathbf{a} = (a_1, a_2, a_3)$ ,  $\mathbf{b} = (b_1, b_2, b_3)$ , and  $\mathbf{c} = (c_1, c_2, c_3)$ . Show the box product,  $[\mathbf{a}, \mathbf{b}, \mathbf{c}]$  equals the determinant

$$\begin{vmatrix} c_1 & c_2 & c_3 \\ a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{vmatrix}.$$

9. It is desired to find an equation of a plane containing the two vectors,  $\mathbf{a}$  and  $\mathbf{b}$ . Using Problem 7, show an equation for this plane is

$$\begin{vmatrix} x & y & z \\ a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{vmatrix} = 0$$

That is, the set of all  $(x, y, z)$  such that the above expression equals zero.

10. Using the notion of the box product yielding either plus or minus the volume of the parallelepiped determined by the given three vectors, show that

$$(\mathbf{a} \times \mathbf{b}) \cdot \mathbf{c} = \mathbf{a} \cdot (\mathbf{b} \times \mathbf{c})$$

In other words, the dot and the cross can be switched as long as the order of the vectors remains the same. **Hint:** There are two ways to do this, by the coordinate description of the dot and cross product and by geometric reasoning.

11. Verify directly that the coordinate description of the cross product,  $\mathbf{a} \times \mathbf{b}$  has the property that it is perpendicular to both  $\mathbf{a}$  and  $\mathbf{b}$ . Then show by direct computation that this coordinate description satisfies

$$\begin{aligned} |\mathbf{a} \times \mathbf{b}|^2 &= |\mathbf{a}|^2 |\mathbf{b}|^2 - (\mathbf{a} \cdot \mathbf{b})^2 \\ &= |\mathbf{a}|^2 |\mathbf{b}|^2 (1 - \cos^2(\theta)) \end{aligned}$$

where  $\theta$  is the angle included between the two vectors. Explain why  $|\mathbf{a} \times \mathbf{b}|$  has the correct magnitude. All that is missing is the material about the right hand rule. Verify directly from the coordinate description of the cross product that the right thing happens with regards to the vectors  $\mathbf{i}, \mathbf{j}, \mathbf{k}$ . Next verify that the distributive law holds for the coordinate description of the cross product. This gives another way to approach the cross product. First define it in terms of coordinates and then get the geometric properties from this.

## 13.6 Vector Identities And Notation

There are two special symbols,  $\delta_{ij}$  and  $\varepsilon_{ijk}$  which are very useful in dealing with vector identities. To begin with, here is the definition of these symbols.

**Definition 13.6.1** The symbol,  $\delta_{ij}$ , called the Kronecker delta symbol is defined as follows.

$$\delta_{ij} \equiv \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}.$$

With the Kronecker symbol,  $i$  and  $j$  can equal any integer in  $\{1, 2, \dots, n\}$  for any  $n \in \mathbb{N}$ .

**Definition 13.6.2** For  $i, j$ , and  $k$  integers in the set,  $\{1, 2, 3\}$ ,  $\varepsilon_{ijk}$  is defined as follows.

$$\varepsilon_{ijk} \equiv \begin{cases} 1 & \text{if } (i, j, k) = (1, 2, 3), (2, 3, 1), \text{ or } (3, 1, 2) \\ -1 & \text{if } (i, j, k) = (2, 1, 3), (1, 3, 2), \text{ or } (3, 2, 1) \\ 0 & \text{if there are any repeated integers} \end{cases}.$$

The subscripts  $ijk$  and  $ij$  in the above are called indices. A single one is called an index. This symbol,  $\varepsilon_{ijk}$  is also called the permutation symbol.

The way to think of  $\varepsilon_{ijk}$  is that  $\varepsilon_{123} = 1$  and if you switch any two of the numbers in the list  $i, j, k$ , it changes the sign. Thus  $\varepsilon_{ijk} = -\varepsilon_{jik}$  and  $\varepsilon_{ijk} = -\varepsilon_{kji}$  etc. You should check that this rule reduces to the above definition. For example, it immediately implies that if there is a repeated index, the answer is zero. This follows because  $\varepsilon_{iij} = -\varepsilon_{iij}$  and so  $\varepsilon_{iij} = 0$ .

It is useful to use the Einstein summation convention when dealing with these symbols. Simply stated, the convention is that you sum over the repeated index. Thus  $a_i b_i$  means  $\sum_i a_i b_i$ . Also,  $\delta_{ij} x_j$  means  $\sum_j \delta_{ij} x_j = x_i$ . When you use this convention, there is one very important thing to never forget. It is this: Never have an index be repeated more than once. Thus  $a_i b_i$  is all right but  $a_{ii} b_i$  is not. The reason for this is that you end up getting confused about what is meant. If you want to write  $\sum_i a_i b_i c_i$  it is best to simply use the summation notation. There is a very important reduction identity connecting these two symbols.

**Lemma 13.6.3** *The following holds.*

$$\varepsilon_{ijk} \varepsilon_{irs} = (\delta_{jr} \delta_{ks} - \delta_{kr} \delta_{js}).$$

**Proof:** If  $\{j, k\} \neq \{r, s\}$  then every term in the sum on the left must have either  $\varepsilon_{ijk}$  or  $\varepsilon_{irs}$  contains a repeated index. Therefore, the left side equals zero. The right side also equals zero in this case. To see this, note that if the two sets are not equal, then there is one of the indices in one of the sets which is not in the other set. For example, it could be that  $j$  is not equal to either  $r$  or  $s$ . Then the right side equals zero.

Therefore, it can be assumed  $\{j, k\} = \{r, s\}$ . If  $i = r$  and  $j = s$  for  $s \neq r$ , then there is exactly one term in the sum on the left and it equals 1. The right also reduces to 1 in this case. If  $i = s$  and  $j = r$ , there is exactly one term in the sum on the left which is nonzero and it must equal -1. The right side also reduces to -1 in this case. If there is a repeated index in  $\{j, k\}$ , then every term in the sum on the left equals zero. The right also reduces to zero in this case because then  $j = k = r = s$  and so the right side becomes  $(1)(1) - (-1)(-1) = 0$ .

**Proposition 13.6.4** *Let  $\mathbf{u}, \mathbf{v}$  be vectors in  $\mathbb{R}^n$  where the Cartesian coordinates of  $\mathbf{u}$  are  $(u_1, \dots, u_n)$  and the Cartesian coordinates of  $\mathbf{v}$  are  $(v_1, \dots, v_n)$ . Then  $\mathbf{u} \cdot \mathbf{v} = u_i v_i$ . If  $\mathbf{u}, \mathbf{v}$  are vectors in  $\mathbb{R}^3$ , then*

$$(\mathbf{u} \times \mathbf{v})_i = \varepsilon_{ijk} u_j v_k.$$

Also,  $\delta_{ik} a_k = a_i$ .

**Proof:** The first claim is obvious from the definition of the dot product. The second is verified by simply checking it works. For example,

$$\mathbf{u} \times \mathbf{v} \equiv \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ u_1 & u_2 & u_3 \\ v_1 & v_2 & v_3 \end{vmatrix}$$

and so

$$(\mathbf{u} \times \mathbf{v})_1 = (u_2 v_3 - u_3 v_2).$$

From the above formula in the proposition,

$$\varepsilon_{1jk} u_j v_k \equiv u_2 v_3 - u_3 v_2,$$

the same thing. The cases for  $(\mathbf{u} \times \mathbf{v})_2$  and  $(\mathbf{u} \times \mathbf{v})_3$  are verified similarly. The last claim follows directly from the definition.

With this notation, you can easily discover vector identities and simplify expressions which involve the cross product.

**Example 13.6.5** Discover a formula which simplifies  $(\mathbf{u} \times \mathbf{v}) \times \mathbf{w}$ .

From the above reduction formula,

$$\begin{aligned}
 ((\mathbf{u} \times \mathbf{v}) \times \mathbf{w})_i &= \varepsilon_{ijk} (\mathbf{u} \times \mathbf{v})_j w_k \\
 &= \varepsilon_{ijk} \varepsilon_{jrs} u_r v_s w_k \\
 &= -\varepsilon_{jik} \varepsilon_{jrs} u_r v_s w_k \\
 &= -(\delta_{ir} \delta_{ks} - \delta_{is} \delta_{kr}) u_r v_s w_k \\
 &= -(u_i v_k w_k - u_k v_i w_k) \\
 &= \mathbf{u} \cdot \mathbf{w} v_i - \mathbf{v} \cdot \mathbf{w} u_i \\
 &= ((\mathbf{u} \cdot \mathbf{w}) \mathbf{v} - (\mathbf{v} \cdot \mathbf{w}) \mathbf{u})_i.
 \end{aligned}$$

Since this holds for all  $i$ , it follows that

$$(\mathbf{u} \times \mathbf{v}) \times \mathbf{w} = (\mathbf{u} \cdot \mathbf{w}) \mathbf{v} - (\mathbf{v} \cdot \mathbf{w}) \mathbf{u}.$$

This is good notation and it will be used in the rest of the book whenever convenient. Actually, this notation is a special case of something more elaborate in which the level of the indices is also important, but there is time for this more general notion later. You will see it in advanced books on mechanics in physics and engineering. It also occurs in the subject of differential geometry.

## 13.7 Exercises

1. Discover a vector identity for  $\mathbf{u} \times (\mathbf{v} \times \mathbf{w})$ .
2. Discover a vector identity for  $(\mathbf{u} \times \mathbf{v}) \cdot (\mathbf{z} \times \mathbf{w})$ .
3. Discover a vector identity for  $(\mathbf{u} \times \mathbf{v}) \times (\mathbf{z} \times \mathbf{w})$  in terms of box products.
4. Simplify  $(\mathbf{u} \times \mathbf{v}) \cdot (\mathbf{v} \times \mathbf{w}) \times (\mathbf{w} \times \mathbf{z})$ .
5. Simplify  $|\mathbf{u} \times \mathbf{v}|^2 + (\mathbf{u} \times \mathbf{v})^2 - |\mathbf{u}|^2 |\mathbf{v}|^2$ .
6. Prove that  $\varepsilon_{ijk} \varepsilon_{ijr} = 2\delta_{kr}$ .

# Functions

Vector valued functions have values in  $\mathbb{R}^p$  where  $p$  is an integer at least as large as 1. Here is a simple example which is obviously of interest.

**Example 14.0.1** *A rocket is launched from the rotating earth. You could define a function having values in  $\mathbb{R}^3$  as  $(r(t), \theta(t), \phi(t))$  where  $r(t)$  is the distance of the center of mass of the rocket from the center of the earth,  $\theta(t)$  is the longitude, and  $\phi(t)$  is the latitude of the rocket.*

**Example 14.0.2** *Let  $\mathbf{f}(x, y) = (\sin xy, y^3 + x, x^4)$ . Then  $\mathbf{f}$  is a function defined on  $\mathbb{R}^2$  which has values in  $\mathbb{R}^3$ . For example,  $\mathbf{f}(1, 2) = (\sin 2, 9, 16)$ .*

As before,  $D(\mathbf{f})$  denotes the domain of the function,  $\mathbf{f}$  which is written in bold face because it will possibly have values in  $\mathbb{R}^p$ . When  $D(\mathbf{f})$  is not specified, it will be understood that the domain of  $\mathbf{f}$  consists of those things for which  $\mathbf{f}$  makes sense.

**Example 14.0.3** *Let  $\mathbf{f}(x, y, z) = (\frac{x+y}{z}, \sqrt{1-x^2}, y)$ . Then  $D(\mathbf{f})$  would consist of the set of all  $(x, y, z)$  such that  $|x| \leq 1$  and  $z \neq 0$ .*

There are many ways to make new functions from old ones.

**Definition 14.0.4** *Let  $\mathbf{f}, \mathbf{g}$  be functions with values in  $\mathbb{R}^p$ . Let  $a, b$  be elements of  $\mathbb{R}$  (scalars). Then  $a\mathbf{f} + b\mathbf{g}$  is the name of a function whose domain is  $D(\mathbf{f}) \cap D(\mathbf{g})$  which is defined as*

$$(a\mathbf{f} + b\mathbf{g})(\mathbf{x}) = a\mathbf{f}(\mathbf{x}) + b\mathbf{g}(\mathbf{x}).$$

$\mathbf{f} \cdot \mathbf{g}$  or  $(\mathbf{f}, \mathbf{g})$  is the name of a function whose domain is  $D(\mathbf{f}) \cap D(\mathbf{g})$  which is defined as

$$(\mathbf{f}, \mathbf{g})(\mathbf{x}) \equiv \mathbf{f} \cdot \mathbf{g}(\mathbf{x}) \equiv \mathbf{f}(\mathbf{x}) \cdot \mathbf{g}(\mathbf{x}).$$

If  $\mathbf{f}$  and  $\mathbf{g}$  have values in  $\mathbb{R}^3$ , define a new function,  $\mathbf{f} \times \mathbf{g}$  by

$$\mathbf{f} \times \mathbf{g}(t) \equiv \mathbf{f}(t) \times \mathbf{g}(t).$$

If  $\mathbf{f} : D(\mathbf{f}) \rightarrow X$  and  $\mathbf{g} : X \rightarrow Y$ , then  $\mathbf{g} \circ \mathbf{f}$  is the name of a function whose domain is

$$\{\mathbf{x} \in D(\mathbf{f}) : \mathbf{f}(\mathbf{x}) \in D(\mathbf{g})\}$$

which is defined as

$$\mathbf{g} \circ \mathbf{f}(\mathbf{x}) \equiv \mathbf{g}(\mathbf{f}(\mathbf{x})).$$

This is called the composition of the two functions.

You should note that  $\mathbf{f}(\mathbf{x})$  is not a function. It is the value of the function at the point,  $\mathbf{x}$ . The name of the function is  $\mathbf{f}$ . Nevertheless, people often write  $\mathbf{f}(\mathbf{x})$  to denote a function and it doesn't cause too many problems in beginning courses. When this is done, the variable,  $\mathbf{x}$  should be considered as a generic variable free to be anything in  $D(\mathbf{f})$ . I will use this slightly sloppy abuse of notation whenever convenient.

**Example 14.0.5** Let  $\mathbf{f}(t) \equiv (t, 1+t, 2)$  and  $\mathbf{g}(t) \equiv (t^2, t, t)$ . Then  $\mathbf{f} \cdot \mathbf{g}$  is the name of the function satisfying

$$\mathbf{f} \cdot \mathbf{g}(t) = \mathbf{f}(t) \cdot \mathbf{g}(t) = t^3 + t + t^2 + 2t = t^3 + t^2 + 3t$$

Note that in this case it was assumed the domains of the functions consisted of all of  $\mathbb{R}$  because this was the set on which the two both made sense. Also note that  $\mathbf{f}$  and  $\mathbf{g}$  map  $\mathbb{R}$  into  $\mathbb{R}^3$  but  $\mathbf{f} \cdot \mathbf{g}$  maps  $\mathbb{R}$  into  $\mathbb{R}$ .

**Example 14.0.6** Suppose  $\mathbf{f}(t) = (2t, 1+t^2)$  and  $g: \mathbb{R}^2 \rightarrow \mathbb{R}$  is given by  $g(x, y) \equiv x + y$ . Then  $g \circ \mathbf{f}: \mathbb{R} \rightarrow \mathbb{R}$  and

$$g \circ \mathbf{f}(t) = g(\mathbf{f}(t)) = g(2t, 1+t^2) = 1 + 2t + t^2.$$

## 14.1 Exercises

1. Let  $\mathbf{f}(t) = \left(t, t^2 + 1, \frac{t}{t+1}\right)$  and let  $\mathbf{g}(t) = \left(t + 1, 1, \frac{t}{t^2+1}\right)$ . Find  $\mathbf{f} \cdot \mathbf{g}$ .
2. Let  $\mathbf{f}, \mathbf{g}$  be given in the previous problem. Find  $\mathbf{f} \times \mathbf{g}$ .
3. Find  $D(\mathbf{f})$  if  $\mathbf{f}(x, y, z, w) = \left(\frac{xy}{zw}, \sqrt{6 - x^2y^2}\right)$ .
4. Let  $\mathbf{f}(t) = (t, t^2, t^3)$ ,  $\mathbf{g}(t) = (1, t, t^2)$ , and  $\mathbf{h}(t) = (\sin t, t, 1)$ . Find the time rate of change of the volume of the parallelepiped spanned by the vectors  $\mathbf{f}, \mathbf{g}$ , and  $\mathbf{h}$ .

## 14.2 Continuous Functions

What was done earlier for scalar functions is generalized here to include the case of a vector valued function.

**Definition 14.2.1** A function  $\mathbf{f}: D(\mathbf{f}) \subseteq \mathbb{R}^p \rightarrow \mathbb{R}^q$  is continuous at  $\mathbf{x} \in D(\mathbf{f})$  if for each  $\varepsilon > 0$  there exists  $\delta > 0$  such that whenever  $\mathbf{y} \in D(\mathbf{f})$  and

$$|\mathbf{y} - \mathbf{x}| < \delta$$

it follows that

$$|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| < \varepsilon.$$

$\mathbf{f}$  is continuous if it is continuous at every point of  $D(\mathbf{f})$ .

Note the total similarity to the scalar valued case.

### 14.2.1 Sufficient Conditions For Continuity

The next theorem is a fundamental result which will allow us to worry less about the  $\varepsilon$   $\delta$  definition of continuity.

**Theorem 14.2.2** *The following assertions are valid.*

1. *The function,  $a\mathbf{f} + b\mathbf{g}$  is continuous at  $\mathbf{x}$  whenever  $\mathbf{f}, \mathbf{g}$  are continuous at  $\mathbf{x} \in D(\mathbf{f}) \cap D(\mathbf{g})$  and  $a, b \in \mathbb{R}$ .*
2. *If  $\mathbf{f}$  is continuous at  $\mathbf{x}$ ,  $\mathbf{f}(\mathbf{x}) \in D(\mathbf{g}) \subseteq \mathbb{R}^p$ , and  $\mathbf{g}$  is continuous at  $\mathbf{f}(\mathbf{x})$ , then  $\mathbf{g} \circ \mathbf{f}$  is continuous at  $\mathbf{x}$ .*
3. *If  $\mathbf{f} = (f_1, \dots, f_q) : D(\mathbf{f}) \rightarrow \mathbb{R}^q$ , then  $\mathbf{f}$  is continuous if and only if each  $f_k$  is a continuous real valued function.*
4. *The function  $f : \mathbb{R}^p \rightarrow \mathbb{R}$ , given by  $f(\mathbf{x}) = |\mathbf{x}|$  is continuous.*

The proof of this theorem is in the last section of this chapter. Its conclusions are not surprising. For example the first claim says that  $(a\mathbf{f} + b\mathbf{g})(\mathbf{y})$  is close to  $(a\mathbf{f} + b\mathbf{g})(\mathbf{x})$  when  $\mathbf{y}$  is close to  $\mathbf{x}$  provided the same can be said about  $\mathbf{f}$  and  $\mathbf{g}$ . For the second claim, if  $\mathbf{y}$  is close to  $\mathbf{x}$ ,  $\mathbf{f}(\mathbf{x})$  is close to  $\mathbf{f}(\mathbf{y})$  and so by continuity of  $\mathbf{g}$  at  $\mathbf{f}(\mathbf{x})$ ,  $\mathbf{g}(\mathbf{f}(\mathbf{y}))$  is close to  $\mathbf{g}(\mathbf{f}(\mathbf{x}))$ . To see the third claim is likely, note that closeness in  $\mathbb{R}^p$  is the same as closeness in each coordinate. The fourth claim is immediate from the triangle inequality.

For functions defined on  $\mathbb{R}^n$ , there is a notion of polynomial just as there is for functions defined on  $\mathbb{R}$ .

**Definition 14.2.3** *Let  $\alpha$  be an  $n$  dimensional multi-index. This means*

$$\alpha = (\alpha_1, \dots, \alpha_n)$$

*where each  $\alpha_i$  is a natural number or zero. Also, let*

$$|\alpha| \equiv \sum_{i=1}^n |\alpha_i|$$

*The symbol,  $\mathbf{x}^\alpha$ , means*

$$\mathbf{x}^\alpha \equiv x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_n^{\alpha_n}.$$

*An  $n$  dimensional polynomial of degree  $m$  is a function of the form*

$$p(\mathbf{x}) = \sum_{|\alpha| \leq m} d_\alpha \mathbf{x}^\alpha.$$

*where the  $d_\alpha$  are real numbers.*

The above theorem implies that polynomials are all continuous.

## 14.3 Exercises

1. Let  $\mathbf{f}(t) = (t, \sin t)$ . Show  $f$  is continuous at every point  $t$ .
2. Suppose  $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| \leq K|\mathbf{x} - \mathbf{y}|$  where  $K$  is a constant. Show that  $\mathbf{f}$  is everywhere continuous. Functions satisfying such an inequality are called Lipschitz functions.

3. Suppose  $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| \leq K|\mathbf{x} - \mathbf{y}|^\alpha$  where  $K$  is a constant and  $\alpha \in (0, 1)$ . Show that  $\mathbf{f}$  is everywhere continuous.
4. Suppose  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  is given by  $f(\mathbf{x}) = 3x_1x_2 + 2x_3^2$ . Use Theorem 14.2.2 to verify that  $f$  is continuous. **Hint:** You should first verify that the function,  $\pi_k : \mathbb{R}^3 \rightarrow \mathbb{R}$  given by  $\pi_k(\mathbf{x}) = x_k$  is a continuous function.
5. Generalize the previous problem to the case where  $f : \mathbb{R}^q \rightarrow \mathbb{R}$  is a polynomial.
6. State and prove a theorem using Theorem 5.4.1 which involves quotients of functions encountered in the previous problem.

## 14.4 Limits Of A Function

As in the case of scalar valued functions of one variable, a concept closely related to continuity is that of the limit of a function. The notion of limit of a function makes sense at points,  $\mathbf{x}$ , which are limit points of  $D(\mathbf{f})$  and this concept is defined next.

**Definition 14.4.1** Let  $A \subseteq \mathbb{R}^m$  be a set. A point,  $\mathbf{x}$ , is a limit point of  $A$  if  $B(\mathbf{x}, r)$  contains infinitely many points of  $A$  for every  $r > 0$ .

**Definition 14.4.2** Let  $\mathbf{f} : D(\mathbf{f}) \subseteq \mathbb{R}^p \rightarrow \mathbb{R}^q$  be a function and let  $\mathbf{x}$  be a limit point of  $D(\mathbf{f})$ . Then

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L}$$

if and only if the following condition holds. For all  $\varepsilon > 0$  there exists  $\delta > 0$  such that if

$$0 < |\mathbf{y} - \mathbf{x}| < \delta, \text{ and } \mathbf{y} \in D(\mathbf{f})$$

then,

$$|\mathbf{L} - \mathbf{f}(\mathbf{y})| < \varepsilon.$$

**Theorem 14.4.3** If  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L}$  and  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L}_1$ , then  $\mathbf{L} = \mathbf{L}_1$ .

**Proof:** Let  $\varepsilon > 0$  be given. There exists  $\delta > 0$  such that if  $0 < |\mathbf{y} - \mathbf{x}| < \delta$  and  $\mathbf{y} \in D(\mathbf{f})$ , then

$$|\mathbf{f}(\mathbf{y}) - \mathbf{L}| < \varepsilon, \quad |\mathbf{f}(\mathbf{y}) - \mathbf{L}_1| < \varepsilon.$$

Pick such a  $\mathbf{y}$ . There exists one because  $\mathbf{x}$  is a limit point of  $D(\mathbf{f})$ . Then

$$|\mathbf{L} - \mathbf{L}_1| \leq |\mathbf{L} - \mathbf{f}(\mathbf{y})| + |\mathbf{f}(\mathbf{y}) - \mathbf{L}_1| < \varepsilon + \varepsilon = 2\varepsilon.$$

Since  $\varepsilon > 0$  was arbitrary, this shows  $\mathbf{L} = \mathbf{L}_1$ .

As in the case of functions of one variable, one can define what it means for  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} f(\mathbf{x}) = \pm\infty$ .

**Definition 14.4.4** If  $f(\mathbf{x}) \in \mathbb{R}$ ,  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} f(\mathbf{x}) = \infty$  if for every number  $l$ , there exists  $\delta > 0$  such that whenever  $|\mathbf{y} - \mathbf{x}| < \delta$  and  $\mathbf{y} \in D(\mathbf{f})$ , then  $f(\mathbf{x}) > l$ .

The following theorem is just like the one variable version presented earlier.



**Theorem 14.4.5** Suppose  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L}$  and  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{g}(\mathbf{y}) = \mathbf{K}$  where  $\mathbf{K}, \mathbf{L} \in \mathbb{R}^q$ . Then if  $a, b \in \mathbb{R}$ ,

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} (a\mathbf{f}(\mathbf{y}) + b\mathbf{g}(\mathbf{y})) = a\mathbf{L} + b\mathbf{K}, \quad (14.1)$$

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) \cdot \mathbf{g}(\mathbf{y}) = \mathbf{L} \cdot \mathbf{K} \quad (14.2)$$

and if  $g$  is scalar valued with  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} g(\mathbf{y}) = K \neq 0$ ,

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) g(\mathbf{y}) = \mathbf{L}K. \quad (14.3)$$

Also, if  $\mathbf{h}$  is a continuous function defined near  $\mathbf{L}$ , then

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{h} \circ \mathbf{f}(\mathbf{y}) = \mathbf{h}(\mathbf{L}). \quad (14.4)$$

Suppose  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L}$ . If  $|\mathbf{f}(\mathbf{y}) - \mathbf{b}| \leq r$  for all  $\mathbf{y}$  sufficiently close to  $\mathbf{x}$ , then  $|\mathbf{L} - \mathbf{b}| \leq r$  also.

**Proof:** The proof of (14.1) is left for you. It is like a corresponding theorem for continuous functions. Now (14.2) is to be verified. Let  $\varepsilon > 0$  be given. Then by the triangle inequality,

$$\begin{aligned} |\mathbf{f} \cdot \mathbf{g}(\mathbf{y}) - \mathbf{L} \cdot \mathbf{K}| &\leq |\mathbf{f} \cdot \mathbf{g}(\mathbf{y}) - \mathbf{f}(\mathbf{y}) \cdot \mathbf{K}| + |\mathbf{f}(\mathbf{y}) \cdot \mathbf{K} - \mathbf{L} \cdot \mathbf{K}| \\ &\leq |\mathbf{f}(\mathbf{y})| |\mathbf{g}(\mathbf{y}) - \mathbf{K}| + |\mathbf{K}| |\mathbf{f}(\mathbf{y}) - \mathbf{L}|. \end{aligned}$$

There exists  $\delta_1$  such that if  $0 < |\mathbf{y} - \mathbf{x}| < \delta_1$  and  $\mathbf{y} \in D(\mathbf{f})$ , then

$$|\mathbf{f}(\mathbf{y}) - \mathbf{L}| < 1,$$

and so for such  $\mathbf{y}$ , the triangle inequality implies,  $|\mathbf{f}(\mathbf{y})| < 1 + |\mathbf{L}|$ . Therefore, for  $0 < |\mathbf{y} - \mathbf{x}| < \delta_1$ ,

$$|\mathbf{f} \cdot \mathbf{g}(\mathbf{y}) - \mathbf{L} \cdot \mathbf{K}| \leq (1 + |\mathbf{K}| + |\mathbf{L}|) [|\mathbf{g}(\mathbf{y}) - \mathbf{K}| + |\mathbf{f}(\mathbf{y}) - \mathbf{L}|]. \quad (14.5)$$

Now let  $0 < \delta_2$  be such that if  $\mathbf{y} \in D(\mathbf{f})$  and  $0 < |\mathbf{x} - \mathbf{y}| < \delta_2$ ,

$$|\mathbf{f}(\mathbf{y}) - \mathbf{L}| < \frac{\varepsilon}{2(1 + |\mathbf{K}| + |\mathbf{L}|)}, \quad |\mathbf{g}(\mathbf{y}) - \mathbf{K}| < \frac{\varepsilon}{2(1 + |\mathbf{K}| + |\mathbf{L}|)}.$$

Then letting  $0 < \delta \leq \min(\delta_1, \delta_2)$ , it follows from (14.5) that

$$|\mathbf{f} \cdot \mathbf{g}(\mathbf{y}) - \mathbf{L} \cdot \mathbf{K}| < \varepsilon$$

and this proves (14.2).

The proof of (14.3) is left to you.

Consider (14.4). Since  $\mathbf{h}$  is continuous near  $\mathbf{L}$ , it follows that for  $\varepsilon > 0$  given, there exists  $\eta > 0$  such that if  $|\mathbf{y} - \mathbf{L}| < \eta$ , then

$$|\mathbf{h}(\mathbf{y}) - \mathbf{h}(\mathbf{L})| < \varepsilon$$

Now since  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L}$ , there exists  $\delta > 0$  such that if  $0 < |\mathbf{y} - \mathbf{x}| < \delta$ , then

$$|\mathbf{f}(\mathbf{y}) - \mathbf{L}| < \eta.$$

Therefore, if  $0 < |\mathbf{y} - \mathbf{x}| < \delta$ ,

$$|\mathbf{h}(\mathbf{f}(\mathbf{y})) - \mathbf{h}(\mathbf{L})| < \varepsilon.$$

It only remains to verify the last assertion. Assume  $|\mathbf{f}(\mathbf{y}) - \mathbf{b}| \leq r$ . It is required to show that  $|\mathbf{L} - \mathbf{b}| \leq r$ . If this is not true, then  $|\mathbf{L} - \mathbf{b}| > r$ . Consider  $B(\mathbf{L}, |\mathbf{L} - \mathbf{b}| - r)$ . Since  $\mathbf{L}$  is the limit of  $\mathbf{f}$ , it follows  $\mathbf{f}(\mathbf{y}) \in B(\mathbf{L}, |\mathbf{L} - \mathbf{b}| - r)$  whenever  $\mathbf{y} \in D(\mathbf{f})$  is close enough to  $\mathbf{x}$ . Thus, by the triangle inequality,

$$|\mathbf{f}(\mathbf{y}) - \mathbf{L}| < |\mathbf{L} - \mathbf{b}| - r$$

and so

$$\begin{aligned} r &< |\mathbf{L} - \mathbf{b}| - |\mathbf{f}(\mathbf{y}) - \mathbf{L}| \leq \|\mathbf{b} - \mathbf{L}\| - \|\mathbf{f}(\mathbf{y}) - \mathbf{L}\| \\ &\leq \|\mathbf{b} - \mathbf{f}(\mathbf{y})\|, \end{aligned}$$

a contradiction to the assumption that  $\|\mathbf{b} - \mathbf{f}(\mathbf{y})\| \leq r$ .

**Theorem 14.4.6** For  $\mathbf{f} : D(\mathbf{f}) \rightarrow \mathbb{R}^q$  and  $\mathbf{x} \in D(\mathbf{f})$  a limit point of  $D(\mathbf{f})$ ,  $\mathbf{f}$  is continuous at  $\mathbf{x}$  if and only if

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{f}(\mathbf{x}).$$

**Proof:** First suppose  $\mathbf{f}$  is continuous at  $\mathbf{x}$  a limit point of  $D(\mathbf{f})$ . Then for every  $\varepsilon > 0$  there exists  $\delta > 0$  such that if  $|\mathbf{y} - \mathbf{x}| < \delta$  and  $\mathbf{y} \in D(\mathbf{f})$ , then  $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| < \varepsilon$ . In particular, this holds if  $0 < |\mathbf{x} - \mathbf{y}| < \delta$  and this is just the definition of the limit. Hence  $\mathbf{f}(\mathbf{x}) = \lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y})$ .

Next suppose  $\mathbf{x}$  is a limit point of  $D(\mathbf{f})$  and  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{f}(\mathbf{x})$ . This means that if  $\varepsilon > 0$  there exists  $\delta > 0$  such that for  $0 < |\mathbf{x} - \mathbf{y}| < \delta$  and  $\mathbf{y} \in D(\mathbf{f})$ , it follows  $|\mathbf{f}(\mathbf{y}) - \mathbf{f}(\mathbf{x})| < \varepsilon$ . However, if  $\mathbf{y} = \mathbf{x}$ , then  $|\mathbf{f}(\mathbf{y}) - \mathbf{f}(\mathbf{x})| = |\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x})| = 0$  and so whenever  $\mathbf{y} \in D(\mathbf{f})$  and  $|\mathbf{x} - \mathbf{y}| < \delta$ , it follows  $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| < \varepsilon$ , showing  $\mathbf{f}$  is continuous at  $\mathbf{x}$ .

The following theorem is important.

**Theorem 14.4.7** Suppose  $\mathbf{f} : D(\mathbf{f}) \rightarrow \mathbb{R}^q$ . Then for  $\mathbf{x}$  a limit point of  $D(\mathbf{f})$ ,

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L} \tag{14.6}$$

if and only if

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} f_k(\mathbf{y}) = L_k \tag{14.7}$$

where  $\mathbf{f}(\mathbf{y}) \equiv (f_1(\mathbf{y}), \dots, f_p(\mathbf{y}))$  and  $\mathbf{L} \equiv (L_1, \dots, L_p)$ .

**Proof:** Suppose (14.6). Then letting  $\varepsilon > 0$  be given there exists  $\delta > 0$  such that if  $0 < |\mathbf{y} - \mathbf{x}| < \delta$ , it follows

$$|f_k(\mathbf{y}) - L_k| \leq |\mathbf{f}(\mathbf{y}) - \mathbf{L}| < \varepsilon$$

which verifies (14.7).

Now suppose (14.7) holds. Then letting  $\varepsilon > 0$  be given, there exists  $\delta_k$  such that if  $0 < |\mathbf{y} - \mathbf{x}| < \delta_k$ , then

$$|f_k(\mathbf{y}) - L_k| < \frac{\varepsilon}{\sqrt{p}}.$$

Let  $0 < \delta < \min(\delta_1, \dots, \delta_p)$ . Then if  $0 < |\mathbf{y} - \mathbf{x}| < \delta$ , it follows

$$\begin{aligned} |\mathbf{f}(\mathbf{y}) - \mathbf{L}| &= \left( \sum_{k=1}^p |f_k(\mathbf{y}) - L_k|^2 \right)^{1/2} \\ &< \left( \sum_{k=1}^p \frac{\varepsilon^2}{p} \right)^{1/2} = \varepsilon. \end{aligned}$$

This proves the theorem.

This theorem shows it suffices to consider the components of a vector valued function when computing the limit.

**Example 14.4.8** Find  $\lim_{(x,y) \rightarrow (3,1)} \left( \frac{x^2-9}{x-3}, y \right)$ .

It is clear that  $\lim_{(x,y) \rightarrow (3,1)} \frac{x^2-9}{x-3} = 6$  and  $\lim_{(x,y) \rightarrow (3,1)} y = 1$ . Therefore, this limit equals  $(6, 1)$ .

**Example 14.4.9** Find  $\lim_{(x,y) \rightarrow (0,0)} \frac{xy}{x^2+y^2}$ .

First of all observe the domain of the function is  $\mathbb{R}^2 \setminus \{(0,0)\}$ , every point in  $\mathbb{R}^2$  except the origin. Therefore,  $(0,0)$  is a limit point of the domain of the function so it might make sense to take a limit. However, just as in the case of a function of one variable, the limit may not exist. In fact, this is the case here. To see this, take points on the line  $y = 0$ . At these points, the value of the function equals 0. Now consider points on the line  $y = x$  where the value of the function equals  $1/2$ . Since arbitrarily close to  $(0,0)$  there are points where the function equals  $1/2$  and points where the function has the value 0, it follows there can be no limit. Just take  $\varepsilon = 1/10$  for example. You can't be within  $1/10$  of  $1/2$  and also within  $1/10$  of 0 at the same time.

Note it is necessary to rely on the definition of the limit much more than in the case of a function of one variable and it is the case there are no easy ways to do limit problems for functions of more than one variable. It is what it is and you will not deal with these concepts without agony.

## 14.5 Exercises

- Find the following limits if possible

- $\lim_{(x,y) \rightarrow (0,0)} \frac{x^2-y^2}{x^2+y^2}$
- $\lim_{(x,y) \rightarrow (0,0)} \frac{x(x^2-y^2)}{(x^2+y^2)}$
- $\lim_{(x,y) \rightarrow (0,0)} \frac{(x^2-y^4)^2}{(x^2+y^4)^2}$  **Hint:** Consider along  $y = 0$  and along  $x = y^2$ .
- $\lim_{(x,y) \rightarrow (0,0)} x \sin \left( \frac{1}{x^2+y^2} \right)$
- $\lim_{(x,y) \rightarrow (1,2)} \frac{-2yx^2+8yx+34y+3y^3-18y^2+6x^2-13x-20-xy^2-x^3}{-y^2+4y-5-x^2+2x}$ . **Hint:** It might help to write this in terms of the variables  $(s, t) = (x-1, y-2)$ .

- In the definition of limit, why must  $\mathbf{x}$  be a limit point of  $D(\mathbf{f})$ ? **Hint:** If  $\mathbf{x}$  were not a limit point of  $D(\mathbf{f})$ , show there exists  $\delta > 0$  such that  $B(\mathbf{x}, \delta)$  contains no points of  $D(\mathbf{f})$  other than possibly  $\mathbf{x}$  itself. Argue that 33.3 is a limit and that so is 22 and 7 and 11. In other words the concept is totally worthless.

## 14.6 The Limit Of A Sequence

As in the case of real numbers, one can consider the limit of a sequence of points in  $\mathbb{R}^p$ .

**Definition 14.6.1** A sequence  $\{\mathbf{a}_n\}_{n=1}^{\infty}$  converges to  $\mathbf{a}$ , and write

$$\lim_{n \rightarrow \infty} \mathbf{a}_n = \mathbf{a} \text{ or } \mathbf{a}_n \rightarrow \mathbf{a}$$

if and only if for every  $\varepsilon > 0$  there exists  $n_\varepsilon$  such that whenever  $n \geq n_\varepsilon$ ,

$$|\mathbf{a}_n - \mathbf{a}| < \varepsilon.$$

In words the definition says that given any measure of closeness,  $\varepsilon$ , the terms of the sequence are eventually all this close to  $\mathbf{a}$ . There is absolutely no difference between this and the definition for sequences of numbers other than here bold face is used to indicate  $\mathbf{a}_n$  and  $\mathbf{a}$  are points in  $\mathbb{R}^p$ .

**Theorem 14.6.2** If  $\lim_{n \rightarrow \infty} \mathbf{a}_n = \mathbf{a}$  and  $\lim_{n \rightarrow \infty} \mathbf{a}_n = \mathbf{a}_1$  then  $\mathbf{a}_1 = \mathbf{a}$ .

**Proof:** Suppose  $\mathbf{a}_1 \neq \mathbf{a}$ . Then let  $0 < \varepsilon < |\mathbf{a}_1 - \mathbf{a}|/2$  in the definition of the limit. It follows there exists  $n_\varepsilon$  such that if  $n \geq n_\varepsilon$ , then  $|\mathbf{a}_n - \mathbf{a}| < \varepsilon$  and  $|\mathbf{a}_n - \mathbf{a}_1| < \varepsilon$ . Therefore, for such  $n$ ,

$$\begin{aligned} |\mathbf{a}_1 - \mathbf{a}| &\leq |\mathbf{a}_1 - \mathbf{a}_n| + |\mathbf{a}_n - \mathbf{a}| \\ &< \varepsilon + \varepsilon < |\mathbf{a}_1 - \mathbf{a}|/2 + |\mathbf{a}_1 - \mathbf{a}|/2 = |\mathbf{a}_1 - \mathbf{a}|, \end{aligned}$$

a contradiction.

As in the case of a vector valued function, it suffices to consider the components. This is the content of the next theorem.

**Theorem 14.6.3** Let  $\mathbf{a}_n = (a_1^n, \dots, a_p^n) \in \mathbb{R}^p$ . Then  $\lim_{n \rightarrow \infty} \mathbf{a}_n = \mathbf{a} \equiv (a_1, \dots, a_p)$  if and only if for each  $k = 1, \dots, p$ ,

$$\lim_{n \rightarrow \infty} a_k^n = a_k. \quad (14.8)$$

**Proof:** First suppose  $\lim_{n \rightarrow \infty} \mathbf{a}_n = \mathbf{a}$ . Then given  $\varepsilon > 0$  there exists  $n_\varepsilon$  such that if  $n > n_\varepsilon$ , then

$$|a_k^n - a_k| \leq |\mathbf{a}_n - \mathbf{a}| < \varepsilon$$

which establishes (14.8).

Now suppose (14.8) holds for each  $k$ . Then letting  $\varepsilon > 0$  be given there exist  $n_k$  such that if  $n > n_k$ ,

$$|a_k^n - a_k| < \varepsilon/\sqrt{p}.$$

Therefore, letting  $n_\varepsilon > \max(n_1, \dots, n_p)$ , it follows that for  $n > n_\varepsilon$ ,

$$|\mathbf{a}_n - \mathbf{a}| = \left( \sum_{k=1}^n |a_k^n - a_k|^2 \right)^{1/2} < \left( \sum_{k=1}^n \frac{\varepsilon^2}{p} \right)^{1/2} = \varepsilon,$$

showing that  $\lim_{n \rightarrow \infty} \mathbf{a}_n = \mathbf{a}$ . This proves the theorem.

**Example 14.6.4** Let  $\mathbf{a}_n = \left( \frac{1}{n^2+1}, \frac{1}{n} \sin(n), \frac{n^2+3}{3n^2+5n} \right)$ .

It suffices to consider the limits of the components according to the following theorem. Thus the limit is  $(0, 0, 1/3)$ .

**Theorem 14.6.5** Suppose  $\{\mathbf{a}_n\}$  and  $\{\mathbf{b}_n\}$  are sequences and that

$$\lim_{n \rightarrow \infty} \mathbf{a}_n = \mathbf{a} \text{ and } \lim_{n \rightarrow \infty} \mathbf{b}_n = \mathbf{b}.$$

Also suppose  $x$  and  $y$  are real numbers. Then

$$\lim_{n \rightarrow \infty} x\mathbf{a}_n + y\mathbf{b}_n = x\mathbf{a} + y\mathbf{b} \quad (14.9)$$

$$\lim_{n \rightarrow \infty} \mathbf{a}_n \cdot \mathbf{b}_n = \mathbf{a} \cdot \mathbf{b} \quad (14.10)$$

If  $b_n \in \mathbb{R}$ , then

$$\mathbf{a}_n b_n \rightarrow \mathbf{a}b.$$

**Proof:** The first of these claims is left for you to do. To do the second, let  $\varepsilon > 0$  be given and choose  $n_1$  such that if  $n \geq n_1$  then

$$|\mathbf{a}_n - \mathbf{a}| < 1.$$

Then for such  $n$ , the triangle inequality and Cauchy Schwarz inequality imply

$$\begin{aligned} |\mathbf{a}_n \cdot \mathbf{b}_n - \mathbf{a} \cdot \mathbf{b}| &\leq |\mathbf{a}_n \cdot \mathbf{b}_n - \mathbf{a}_n \cdot \mathbf{b}| + |\mathbf{a}_n \cdot \mathbf{b} - \mathbf{a} \cdot \mathbf{b}| \\ &\leq |\mathbf{a}_n| |\mathbf{b}_n - \mathbf{b}| + |\mathbf{b}| |\mathbf{a}_n - \mathbf{a}| \\ &\leq (|\mathbf{a}| + 1) |\mathbf{b}_n - \mathbf{b}| + |\mathbf{b}| |\mathbf{a}_n - \mathbf{a}|. \end{aligned}$$

Now let  $n_2$  be large enough that for  $n \geq n_2$ ,

$$|\mathbf{b}_n - \mathbf{b}| < \frac{\varepsilon}{2(|\mathbf{a}| + 1)}, \text{ and } |\mathbf{a}_n - \mathbf{a}| < \frac{\varepsilon}{2(|\mathbf{b}| + 1)}.$$

Such a number exists because of the definition of limit. Therefore, let

$$n_\varepsilon > \max(n_1, n_2).$$

For  $n \geq n_\varepsilon$ ,

$$\begin{aligned} |\mathbf{a}_n \cdot \mathbf{b}_n - \mathbf{a} \cdot \mathbf{b}| &\leq (|\mathbf{a}| + 1) |\mathbf{b}_n - \mathbf{b}| + |\mathbf{b}| |\mathbf{a}_n - \mathbf{a}| \\ &< (|\mathbf{a}| + 1) \frac{\varepsilon}{2(|\mathbf{a}| + 1)} + |\mathbf{b}| \frac{\varepsilon}{2(|\mathbf{b}| + 1)} \leq \varepsilon. \end{aligned}$$

This proves (14.9). The proof of (14.10) is entirely similar and is left for you.

### 14.6.1 Sequences And Completeness

Recall the definition of a Cauchy sequence.

**Definition 14.6.6**  $\{\mathbf{a}_n\}$  is a Cauchy sequence if for all  $\varepsilon > 0$ , there exists  $n_\varepsilon$  such that whenever  $n, m \geq n_\varepsilon$ ,

$$|\mathbf{a}_n - \mathbf{a}_m| < \varepsilon.$$

A sequence is Cauchy means the terms are “bunching up to each other” as  $m, n$  get large.

**Theorem 14.6.7** Let  $\{\mathbf{a}_n\}_{n=1}^\infty$  be a Cauchy sequence in  $\mathbb{R}^p$ . Then there exists  $\mathbf{a} \in \mathbb{R}^p$  such that  $\mathbf{a}_n \rightarrow \mathbf{a}$ .

**Proof:** Let  $\mathbf{a}_n = (a_1^n, \dots, a_p^n)$ . Then

$$|a_k^n - a_k^m| \leq |\mathbf{a}_n - \mathbf{a}_m|$$

which shows for each  $k = 1, \dots, p$ , it follows  $\{a_k^n\}_{n=1}^\infty$  is a Cauchy sequence. By completeness of  $\mathbb{R}$ , it follows there exists  $a_k$  such that  $\lim_{n \rightarrow \infty} a_k^n = a_k$ . Letting  $\mathbf{a} = (a_1, \dots, a_p)$ , it follows from Theorem 14.6.3 that

$$\lim_{n \rightarrow \infty} \mathbf{a}_n = \mathbf{a}.$$

This proves the theorem.

**Theorem 14.6.8** *The set of terms in a Cauchy sequence in  $\mathbb{R}^p$  is bounded in the sense that for all  $n$ ,  $|\mathbf{a}_n| < M$  for some  $M < \infty$ .*

**Proof:** Let  $\varepsilon = 1$  in the definition of a Cauchy sequence and let  $n > n_1$ . Then from the definition,

$$|\mathbf{a}_n - \mathbf{a}_{n_1}| < 1.$$

It follows that for all  $n > n_1$ ,

$$|\mathbf{a}_n| < 1 + |\mathbf{a}_{n_1}|.$$

Therefore, for all  $n$ ,

$$|\mathbf{a}_n| \leq 1 + |\mathbf{a}_{n_1}| + \sum_{k=1}^{n_1} |\mathbf{a}_k|.$$

This proves the theorem.

**Theorem 14.6.9** *If a sequence  $\{\mathbf{a}_n\}$  in  $\mathbb{R}^p$  converges, then the sequence is a Cauchy sequence.*

**Proof:** Let  $\varepsilon > 0$  be given and suppose  $\mathbf{a}_n \rightarrow \mathbf{a}$ . Then from the definition of convergence, there exists  $n_\varepsilon$  such that if  $n > n_\varepsilon$ , it follows that

$$|\mathbf{a}_n - \mathbf{a}| < \frac{\varepsilon}{2}$$

Therefore, if  $m, n \geq n_\varepsilon + 1$ , it follows that

$$|\mathbf{a}_n - \mathbf{a}_m| \leq |\mathbf{a}_n - \mathbf{a}| + |\mathbf{a} - \mathbf{a}_m| < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon$$

showing that, since  $\varepsilon > 0$  is arbitrary,  $\{\mathbf{a}_n\}$  is a Cauchy sequence.

## 14.6.2 Continuity And The Limit Of A Sequence

Just as in the case of a function of one variable, there is a very useful way of thinking of continuity in terms of limits of sequences found in the following theorem. In words, it says a function is continuous if it takes convergent sequences to convergent sequences whenever possible.

**Theorem 14.6.10** *A function  $\mathbf{f} : D(\mathbf{f}) \rightarrow \mathbb{R}^q$  is continuous at  $\mathbf{x} \in D(\mathbf{f})$  if and only if, whenever  $\mathbf{x}_n \rightarrow \mathbf{x}$  with  $\mathbf{x}_n \in D(\mathbf{f})$ , it follows  $\mathbf{f}(\mathbf{x}_n) \rightarrow \mathbf{f}(\mathbf{x})$ .*

**Proof:** Suppose first that  $\mathbf{f}$  is continuous at  $\mathbf{x}$  and let  $\mathbf{x}_n \rightarrow \mathbf{x}$ . Let  $\varepsilon > 0$  be given. By continuity, there exists  $\delta > 0$  such that if  $|\mathbf{y} - \mathbf{x}| < \delta$ , then  $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| < \varepsilon$ . However, there exists  $n_\delta$  such that if  $n \geq n_\delta$ , then  $|\mathbf{x}_n - \mathbf{x}| < \delta$  and so for all  $n$  this large,

$$|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}_n)| < \varepsilon$$

which shows  $\mathbf{f}(\mathbf{x}_n) \rightarrow \mathbf{f}(\mathbf{x})$ .

Now suppose the condition about taking convergent sequences to convergent sequences holds at  $\mathbf{x}$ . Suppose  $\mathbf{f}$  fails to be continuous at  $\mathbf{x}$ . Then there exists  $\varepsilon > 0$  and  $\mathbf{x}_n \in D(\mathbf{f})$  such that  $|\mathbf{x} - \mathbf{x}_n| < \frac{1}{n}$ , yet

$$|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}_n)| \geq \varepsilon.$$

But this is clearly a contradiction because, although  $\mathbf{x}_n \rightarrow \mathbf{x}$ ,  $\mathbf{f}(\mathbf{x}_n)$  fails to converge to  $\mathbf{f}(\mathbf{x})$ . It follows  $\mathbf{f}$  must be continuous after all. This proves the theorem.

## 14.7 Properties Of Continuous Functions

Functions of  $p$  variables have many of the same properties as functions of one variable. First there is a version of the extreme value theorem generalizing the one dimensional case.

**Theorem 14.7.1** *Let  $C$  be closed and bounded and let  $f : C \rightarrow \mathbb{R}$  be continuous. Then  $f$  achieves its maximum and its minimum on  $C$ . This means there exist,  $\mathbf{x}_1, \mathbf{x}_2 \in C$  such that for all  $\mathbf{x} \in C$ ,*

$$f(\mathbf{x}_1) \leq f(\mathbf{x}) \leq f(\mathbf{x}_2).$$

There is also the long technical theorem about sums and products of continuous functions. These theorems are proved in the next section.

**Theorem 14.7.2** *The following assertions are valid*

1. *The function,  $a\mathbf{f} + b\mathbf{g}$  is continuous at  $\mathbf{x}$  when  $\mathbf{f}, \mathbf{g}$  are continuous at  $\mathbf{x} \in D(\mathbf{f}) \cap D(\mathbf{g})$  and  $a, b \in \mathbb{R}$ .*
2. *If  $f$  and  $g$  are each real valued functions continuous at  $\mathbf{x}$ , then  $fg$  is continuous at  $\mathbf{x}$ . If, in addition to this,  $g(\mathbf{x}) \neq 0$ , then  $f/g$  is continuous at  $\mathbf{x}$ .*
3. *If  $\mathbf{f}$  is continuous at  $\mathbf{x}$ ,  $\mathbf{f}(\mathbf{x}) \in D(\mathbf{g}) \subseteq \mathbb{R}^p$ , and  $\mathbf{g}$  is continuous at  $\mathbf{f}(\mathbf{x})$ , then  $\mathbf{g} \circ \mathbf{f}$  is continuous at  $\mathbf{x}$ .*
4. *If  $\mathbf{f} = (f_1, \dots, f_q) : D(\mathbf{f}) \rightarrow \mathbb{R}^q$ , then  $\mathbf{f}$  is continuous if and only if each  $f_k$  is a continuous real valued function.*
5. *The function  $f : \mathbb{R}^p \rightarrow \mathbb{R}$ , given by  $f(\mathbf{x}) = |\mathbf{x}|$  is continuous.*

## 14.8 Exercises

1.  $\mathbf{f} : D \subseteq \mathbb{R}^p \rightarrow \mathbb{R}^q$  is Lipschitz continuous or just Lipschitz for short if there exists a constant,  $K$  such that

$$|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| \leq K |\mathbf{x} - \mathbf{y}|$$

for all  $\mathbf{x}, \mathbf{y} \in D$ . Show every Lipschitz function is uniformly continuous which means that given  $\varepsilon > 0$  there exists  $\delta > 0$  independent of  $\mathbf{x}$  such that if  $|\mathbf{x} - \mathbf{y}| < \delta$ , then  $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| < \varepsilon$ .

2. If  $\mathbf{f}$  is uniformly continuous, does it follow that  $|\mathbf{f}|$  is also uniformly continuous? If  $|\mathbf{f}|$  is uniformly continuous does it follow that  $\mathbf{f}$  is uniformly continuous? Answer the same questions with “uniformly continuous” replaced with “continuous”. Explain why.

## 14.9 Some Advanced Calculus

This section contains the proofs of the theorems which were just stated without proof.

**Theorem 14.9.1** *The following assertions are valid*

1. *The function,  $a\mathbf{f} + b\mathbf{g}$  is continuous at  $\mathbf{x}$  when  $\mathbf{f}, \mathbf{g}$  are continuous at  $\mathbf{x} \in D(\mathbf{f}) \cap D(\mathbf{g})$  and  $a, b \in \mathbb{R}$ .*
2. *If  $f$  and  $g$  are each real valued functions continuous at  $\mathbf{x}$ , then  $fg$  is continuous at  $\mathbf{x}$ . If, in addition to this,  $g(\mathbf{x}) \neq 0$ , then  $f/g$  is continuous at  $\mathbf{x}$ .*
3. *If  $\mathbf{f}$  is continuous at  $\mathbf{x}$ ,  $\mathbf{f}(\mathbf{x}) \in D(\mathbf{g}) \subseteq \mathbb{R}^p$ , and  $\mathbf{g}$  is continuous at  $\mathbf{f}(\mathbf{x})$ , then  $\mathbf{g} \circ \mathbf{f}$  is continuous at  $\mathbf{x}$ .*
4. *If  $\mathbf{f} = (f_1, \dots, f_q) : D(\mathbf{f}) \rightarrow \mathbb{R}^q$ , then  $\mathbf{f}$  is continuous if and only if each  $f_k$  is a continuous real valued function.*
5. *The function  $f : \mathbb{R}^p \rightarrow \mathbb{R}$ , given by  $f(\mathbf{x}) = |\mathbf{x}|$  is continuous.*

**Proof:** Begin with 1.) Let  $\varepsilon > 0$  be given. By assumption, there exist  $\delta_1 > 0$  such that whenever  $|\mathbf{x} - \mathbf{y}| < \delta_1$ , it follows  $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| < \frac{\varepsilon}{2(|a|+|b|+1)}$  and there exists  $\delta_2 > 0$  such that whenever  $|\mathbf{x} - \mathbf{y}| < \delta_2$ , it follows that  $|\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y})| < \frac{\varepsilon}{2(|a|+|b|+1)}$ . Then let  $0 < \delta \leq \min(\delta_1, \delta_2)$ . If  $|\mathbf{x} - \mathbf{y}| < \delta$ , then everything happens at once. Therefore, using the triangle inequality

$$\begin{aligned} & |a\mathbf{f}(\mathbf{x}) + b\mathbf{f}(\mathbf{x}) - (a\mathbf{g}(\mathbf{y}) + b\mathbf{g}(\mathbf{y}))| \\ & \leq |a| |\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| + |b| |\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y})| \\ & < |a| \left( \frac{\varepsilon}{2(|a|+|b|+1)} \right) + |b| \left( \frac{\varepsilon}{2(|a|+|b|+1)} \right) < \varepsilon. \end{aligned}$$

Now begin on 2.) There exists  $\delta_1 > 0$  such that if  $|\mathbf{y} - \mathbf{x}| < \delta_1$ , then  $|f(\mathbf{x}) - f(\mathbf{y})| < 1$ . Therefore, for such  $\mathbf{y}$ ,

$$|f(\mathbf{y})| < 1 + |f(\mathbf{x})|.$$

It follows that for such  $\mathbf{y}$ ,

$$\begin{aligned} |fg(\mathbf{x}) - fg(\mathbf{y})| & \leq |f(\mathbf{x})g(\mathbf{x}) - g(\mathbf{x})f(\mathbf{y})| + |g(\mathbf{x})f(\mathbf{y}) - f(\mathbf{y})g(\mathbf{y})| \\ & \leq |g(\mathbf{x})| |f(\mathbf{x}) - f(\mathbf{y})| + |f(\mathbf{y})| |g(\mathbf{x}) - g(\mathbf{y})| \\ & \leq (1 + |g(\mathbf{x})| + |f(\mathbf{y})|) [|g(\mathbf{x}) - g(\mathbf{y})| + |f(\mathbf{x}) - f(\mathbf{y})|]. \end{aligned}$$

Now let  $\varepsilon > 0$  be given. There exists  $\delta_2$  such that if  $|\mathbf{x} - \mathbf{y}| < \delta_2$ , then

$$|g(\mathbf{x}) - g(\mathbf{y})| < \frac{\varepsilon}{2(1 + |g(\mathbf{x})| + |f(\mathbf{y})|)},$$

and there exists  $\delta_3$  such that if  $|\mathbf{x} - \mathbf{y}| < \delta_3$ , then

$$|f(\mathbf{x}) - f(\mathbf{y})| < \frac{\varepsilon}{2(1 + |g(\mathbf{x})| + |f(\mathbf{y})|)}$$

Now let  $0 < \delta \leq \min(\delta_1, \delta_2, \delta_3)$ . Then if  $|\mathbf{x} - \mathbf{y}| < \delta$ , all the above hold at once and

$$|fg(\mathbf{x}) - fg(\mathbf{y})| \leq$$



$$\begin{aligned}
& (1 + |g(\mathbf{x})| + |f(\mathbf{y})|) [|g(\mathbf{x}) - g(\mathbf{y})| + |f(\mathbf{x}) - f(\mathbf{y})|] \\
& < (1 + |g(\mathbf{x})| + |f(\mathbf{y})|) \left( \frac{\varepsilon}{2(1 + |g(\mathbf{x})| + |f(\mathbf{y})|)} + \frac{\varepsilon}{2(1 + |g(\mathbf{x})| + |f(\mathbf{y})|)} \right) = \varepsilon.
\end{aligned}$$

This proves the first part of 2.) To obtain the second part, let  $\delta_1$  be as described above and let  $\delta_0 > 0$  be such that for  $|\mathbf{x} - \mathbf{y}| < \delta_0$ ,

$$|g(\mathbf{x}) - g(\mathbf{y})| < |g(\mathbf{x})|/2$$

and so by the triangle inequality,

$$-|g(\mathbf{x})|/2 \leq |g(\mathbf{y})| - |g(\mathbf{x})| \leq |g(\mathbf{x})|/2$$

which implies  $|g(\mathbf{y})| \geq |g(\mathbf{x})|/2$ , and  $|g(\mathbf{y})| < 3|g(\mathbf{x})|/2$ .

Then if  $|\mathbf{x} - \mathbf{y}| < \min(\delta_0, \delta_1)$ ,

$$\begin{aligned}
\left| \frac{f(\mathbf{x})}{g(\mathbf{x})} - \frac{f(\mathbf{y})}{g(\mathbf{y})} \right| &= \left| \frac{f(\mathbf{x})g(\mathbf{y}) - f(\mathbf{y})g(\mathbf{x})}{g(\mathbf{x})g(\mathbf{y})} \right| \\
&\leq \frac{|f(\mathbf{x})g(\mathbf{y}) - f(\mathbf{y})g(\mathbf{x})|}{\left(\frac{|g(\mathbf{x})|^2}{2}\right)} \\
&= \frac{2|f(\mathbf{x})g(\mathbf{y}) - f(\mathbf{y})g(\mathbf{x})|}{|g(\mathbf{x})|^2} \\
&\leq \frac{2}{|g(\mathbf{x})|^2} [|f(\mathbf{x})g(\mathbf{y}) - f(\mathbf{y})g(\mathbf{y}) + f(\mathbf{y})g(\mathbf{y}) - f(\mathbf{y})g(\mathbf{x})|] \\
&\leq \frac{2}{|g(\mathbf{x})|^2} [|g(\mathbf{y})||f(\mathbf{x}) - f(\mathbf{y})| + |f(\mathbf{y})||g(\mathbf{y}) - g(\mathbf{x})|] \\
&\leq \frac{2}{|g(\mathbf{x})|^2} \left[ \frac{3}{2}|g(\mathbf{x})||f(\mathbf{x}) - f(\mathbf{y})| + (1 + |f(\mathbf{x})|)|g(\mathbf{y}) - g(\mathbf{x})| \right] \\
&\leq \frac{2}{|g(\mathbf{x})|^2} (1 + 2|f(\mathbf{x})| + 2|g(\mathbf{x})|) [|f(\mathbf{x}) - f(\mathbf{y})| + |g(\mathbf{y}) - g(\mathbf{x})|] \\
&\equiv M [|f(\mathbf{x}) - f(\mathbf{y})| + |g(\mathbf{y}) - g(\mathbf{x})|]
\end{aligned}$$

where

$$M \equiv \frac{2}{|g(\mathbf{x})|^2} (1 + 2|f(\mathbf{x})| + 2|g(\mathbf{x})|)$$

Now let  $\delta_2$  be such that if  $|\mathbf{x} - \mathbf{y}| < \delta_2$ , then

$$|f(\mathbf{x}) - f(\mathbf{y})| < \frac{\varepsilon}{2} M^{-1}$$

and let  $\delta_3$  be such that if  $|\mathbf{x} - \mathbf{y}| < \delta_3$ , then

$$|g(\mathbf{y}) - g(\mathbf{x})| < \frac{\varepsilon}{2} M^{-1}.$$

Then if  $0 < \delta \leq \min(\delta_0, \delta_1, \delta_2, \delta_3)$ , and  $|\mathbf{x} - \mathbf{y}| < \delta$ , everything holds and

$$\left| \frac{f(\mathbf{x})}{g(\mathbf{x})} - \frac{f(\mathbf{y})}{g(\mathbf{y})} \right| \leq M [|f(\mathbf{x}) - f(\mathbf{y})| + |g(\mathbf{y}) - g(\mathbf{x})|]$$

$$< M \left[ \frac{\varepsilon}{2} M^{-1} + \frac{\varepsilon}{2} M^{-1} \right] = \varepsilon.$$

This completes the proof of the second part of 2.) Note that in these proofs no effort is made to find some sort of “best”  $\delta$ . The problem is one which has a yes or a no answer. Either it is or it is not continuous.

Now begin on 3.). If  $\mathbf{f}$  is continuous at  $\mathbf{x}$ ,  $\mathbf{f}(\mathbf{x}) \in D(\mathbf{g}) \subseteq \mathbb{R}^p$ , and  $\mathbf{g}$  is continuous at  $\mathbf{f}(\mathbf{x})$ , then  $\mathbf{g} \circ \mathbf{f}$  is continuous at  $\mathbf{x}$ . Let  $\varepsilon > 0$  be given. Then there exists  $\eta > 0$  such that if  $|\mathbf{y} - \mathbf{f}(\mathbf{x})| < \eta$  and  $\mathbf{y} \in D(\mathbf{g})$ , it follows that  $|\mathbf{g}(\mathbf{y}) - \mathbf{g}(\mathbf{f}(\mathbf{x}))| < \varepsilon$ . It follows from continuity of  $\mathbf{f}$  at  $\mathbf{x}$  that there exists  $\delta > 0$  such that if  $|\mathbf{x} - \mathbf{z}| < \delta$  and  $\mathbf{z} \in D(\mathbf{f})$ , then  $|\mathbf{f}(\mathbf{z}) - \mathbf{f}(\mathbf{x})| < \eta$ . Then if  $|\mathbf{x} - \mathbf{z}| < \delta$  and  $\mathbf{z} \in D(\mathbf{g} \circ \mathbf{f}) \subseteq D(\mathbf{f})$ , all the above hold and so

$$|\mathbf{g}(\mathbf{f}(\mathbf{z})) - \mathbf{g}(\mathbf{f}(\mathbf{x}))| < \varepsilon.$$

This proves part 3.)

Part 4.) says: If  $\mathbf{f} = (f_1, \dots, f_q) : D(\mathbf{f}) \rightarrow \mathbb{R}^q$ , then  $\mathbf{f}$  is continuous if and only if each  $f_k$  is a continuous real valued function. Then

$$\begin{aligned} |f_k(\mathbf{x}) - f_k(\mathbf{y})| &\leq |\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| \\ &\equiv \left( \sum_{i=1}^q |f_i(\mathbf{x}) - f_i(\mathbf{y})|^2 \right)^{1/2} \\ &\leq \sum_{i=1}^q |f_i(\mathbf{x}) - f_i(\mathbf{y})|. \end{aligned} \quad (14.11)$$

Suppose first that  $\mathbf{f}$  is continuous at  $\mathbf{x}$ . Then there exists  $\delta > 0$  such that if  $|\mathbf{x} - \mathbf{y}| < \delta$ , then  $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| < \varepsilon$ . The first part of the above inequality then shows that for each  $k = 1, \dots, q$ ,  $|f_k(\mathbf{x}) - f_k(\mathbf{y})| < \varepsilon$ . This shows the only if part. Now suppose each function,  $f_k$  is continuous. Then if  $\varepsilon > 0$  is given, there exists  $\delta_k > 0$  such that whenever  $|\mathbf{x} - \mathbf{y}| < \delta_k$

$$|f_k(\mathbf{x}) - f_k(\mathbf{y})| < \varepsilon/q.$$

Now let  $0 < \delta \leq \min(\delta_1, \dots, \delta_q)$ . For  $|\mathbf{x} - \mathbf{y}| < \delta$ , the above inequality holds for all  $k$  and so the last part of (14.11) implies

$$\begin{aligned} |\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| &\leq \sum_{i=1}^q |f_i(\mathbf{x}) - f_i(\mathbf{y})| \\ &< \sum_{i=1}^q \frac{\varepsilon}{q} = \varepsilon. \end{aligned}$$

This proves part 4.)

To verify part 5.), let  $\varepsilon > 0$  be given and let  $\delta = \varepsilon$ . Then if  $|\mathbf{x} - \mathbf{y}| < \delta$ , the triangle inequality implies

$$\begin{aligned} |f(\mathbf{x}) - f(\mathbf{y})| &= ||\mathbf{x}| - |\mathbf{y}|| \\ &\leq |\mathbf{x} - \mathbf{y}| < \delta = \varepsilon. \end{aligned}$$

This proves part 5.) and completes the proof of the theorem.

Here is a multidimensional version of the nested interval lemma.

**Lemma 14.9.2** Let  $I_k = \prod_{i=1}^p [a_i^k, b_i^k] \equiv \{\mathbf{x} \in \mathbb{R}^p : x_i \in [a_i^k, b_i^k]\}$  and suppose that for all  $k = 1, 2, \dots$ ,

$$I_k \supseteq I_{k+1}.$$

Then there exists a point,  $\mathbf{c} \in \mathbb{R}^p$  which is an element of every  $I_k$ .

**Proof:** Since  $I_k \supseteq I_{k+1}$ , it follows that for each  $i = 1, \dots, p$ ,  $[a_i^k, b_i^k] \supseteq [a_i^{k+1}, b_i^{k+1}]$ . This implies that for each  $i$ ,

$$a_i^k \leq a_i^{k+1}, \quad b_i^k \geq b_i^{k+1}. \quad (14.12)$$

Consequently, if  $k \leq l$ ,

$$a_i^l \leq b_i^l \leq b_i^k. \quad (14.13)$$

Now define

$$c_i \equiv \sup \{a_i^l : l = 1, 2, \dots\}$$

By the first inequality in (14.12),

$$c_i = \sup \{a_i^l : l = k, k+1, \dots\} \quad (14.14)$$

for each  $k = 1, 2, \dots$ . Therefore, picking any  $k$ , (14.13) shows that  $b_i^k$  is an upper bound for the set,  $\{a_i^l : l = k, k+1, \dots\}$  and so it is at least as large as the least upper bound of this set which is the definition of  $c_i$  given in (14.14). Thus, for each  $i$  and each  $k$ ,

$$a_i^k \leq c_i \leq b_i^k.$$

Defining  $\mathbf{c} \equiv (c_1, \dots, c_p)$ ,  $\mathbf{c} \in I_k$  for all  $k$ . This proves the lemma.

If you don't like the proof, here is another one which follows directly from the one variable nested interval lemma, Lemma 5.7.7 on Page 95.

**Lemma 14.9.3** Let  $I_k = \prod_{i=1}^p [a_i^k, b_i^k] \equiv \{\mathbf{x} \in \mathbb{R}^p : x_i \in [a_i^k, b_i^k]\}$  and suppose that for all  $k = 1, 2, \dots$ ,

$$I_k \supseteq I_{k+1}.$$

Then there exists a point,  $\mathbf{c} \in \mathbb{R}^p$  which is an element of every  $I_k$ .

**Proof:** For each  $i = 1, \dots, p$ ,  $[a_i^k, b_i^k] \supseteq [a_i^{k+1}, b_i^{k+1}]$  and so by the nested interval theorem for one dimensional problems, there exists a point  $c_i \in [a_i^k, b_i^k]$  for all  $k$ . Then letting  $\mathbf{c} \equiv (c_1, \dots, c_p)$  it follows  $\mathbf{c} \in I_k$  for all  $k$ . This proves the lemma.

The following definition is similar to that given earlier. It defines what is meant by a sequentially compact set in  $\mathbb{R}^p$ .

**Definition 14.9.4** A set,  $K \subseteq \mathbb{R}^p$  is sequentially compact if and only if whenever  $\{\mathbf{x}_n\}_{n=1}^\infty$  is a sequence of points in  $K$ , there exists a point,  $\mathbf{x} \in K$  and a subsequence,  $\{\mathbf{x}_{n_k}\}_{k=1}^\infty$  such that  $\mathbf{x}_{n_k} \rightarrow \mathbf{x}$ .

It turns out the sequentially compact sets in  $\mathbb{R}^p$  are exactly those which are closed and bounded. Only half of this result will be needed in this book and this is proved next.

**Theorem 14.9.5** Let  $C \subseteq \mathbb{R}^p$  be closed and bounded. Then  $C$  is sequentially compact.

**Proof:** Let  $\{\mathbf{a}_n\} \subseteq C$ , let  $C \subseteq \prod_{i=1}^p [a_i, b_i]$ , and consider all sets of the form  $\prod_{i=1}^p [c_i, d_i]$  where  $[c_i, d_i]$  equals either  $[a_i, \frac{a_i+b_i}{2}]$  or  $[c_i, d_i] = [\frac{a_i+b_i}{2}, b_i]$ . Thus there are  $2^p$  of these sets because there are two choices for the  $i^{\text{th}}$  slot for  $i = 1, \dots, p$ . Also, if  $\mathbf{x}$  and  $\mathbf{y}$  are two points in one of these sets,

$$|x_i - y_i| \leq 2^{-1} |b_i - a_i|.$$

Therefore, letting  $D_0 = \left(\sum_{i=1}^p |b_i - a_i|^2\right)^{1/2}$ ,

$$\begin{aligned} |\mathbf{x} - \mathbf{y}| &= \left(\sum_{i=1}^p |x_i - y_i|^2\right)^{1/2} \\ &\leq 2^{-1} \left(\sum_{i=1}^p |b_i - a_i|^2\right)^{1/2} \equiv 2^{-1} D_0. \end{aligned}$$

In particular, since  $\mathbf{d} \equiv (d_1, \dots, d_p)$  and  $\mathbf{c} \equiv (c_1, \dots, c_p)$  are two such points,

$$D_1 \equiv \left( \sum_{i=1}^p |d_i - c_i|^2 \right)^{1/2} \leq 2^{-1} D_0$$

Denote by  $\{J_1, \dots, J_{2^p}\}$  these sets determined above. Since the union of these sets equals all of  $I_0$ , it follows

$$C = \bigcup_{k=1}^{2^p} J_k \cap C.$$

Pick  $J_k$  such that  $\mathbf{a}_n$  is contained in  $J_k \cap C$  for infinitely many values of  $n$ . Let  $I_1 \equiv J_k$ . Now do to  $I_1$  what was done to  $I_0$  to obtain  $I_2 \subseteq I_1$  and for any two points,  $\mathbf{x}, \mathbf{y} \in I_2$

$$|\mathbf{x} - \mathbf{y}| \leq 2^{-1} D_1 \leq 2^{-2} D_0,$$

and  $I_2 \cap C$  contains  $\mathbf{a}_n$  for infinitely many values of  $n$ . Continue in this way obtaining sets,  $I_k$  such that  $I_k \supseteq I_{k+1}$  and for any two points in  $I_k$ ,  $\mathbf{x}, \mathbf{y}$ , it follows  $|\mathbf{x} - \mathbf{y}| \leq 2^{-k} D_0$ , and  $I_k \cap C$  contains  $\mathbf{a}_n$  for infinitely many values of  $n$ . By the nested interval lemma, there exists a point,  $\mathbf{c}$  which is contained in each  $I_k$ .

**Claim:**  $\mathbf{c} \in C$ .

**Proof of claim:** Suppose  $\mathbf{c} \notin C$ . Since  $C$  is a closed set, there exists  $r > 0$  such that  $B(\mathbf{c}, r)$  is contained completely in  $\mathbb{R}^p \setminus C$ . In other words,  $B(\mathbf{c}, r)$  contains no points of  $C$ . Let  $k$  be so large that  $D_0 2^{-k} < r$ . Then since  $\mathbf{c} \in I_k$ , and any two points of  $I_k$  are closer than  $D_0 2^{-k}$ ,  $I_k$  must be contained in  $B(\mathbf{c}, r)$  and so has no points of  $C$  in it, contrary to the manner in which the  $I_k$  are defined in which  $I_k$  contains  $\mathbf{a}_n$  for infinitely many values of  $n$ . Therefore,  $\mathbf{c} \in C$  as claimed.

Now pick  $\mathbf{a}_{n_1} \in I_1 \cap \{\mathbf{a}_n\}_{n=1}^\infty$ . Having picked this, let  $\mathbf{a}_{n_2} \in I_2 \cap \{\mathbf{a}_n\}_{n=1}^\infty$  with  $n_2 > n_1$ . Having picked these two, let  $\mathbf{a}_{n_3} \in I_3 \cap \{\mathbf{a}_n\}_{n=1}^\infty$  with  $n_3 > n_2$  and continue this way. The result is a subsequence of  $\{\mathbf{a}_n\}_{n=1}^\infty$  which converges to  $\mathbf{c} \in C$  because any two points in  $I_k$  are within  $D_0 2^{-k}$  of each other. This proves the theorem.

Here is a proof of the extreme value theorem.

**Theorem 14.9.6** *Let  $C$  be closed and bounded and let  $f : C \rightarrow \mathbb{R}$  be continuous. Then  $f$  achieves its maximum and its minimum on  $C$ . This means there exist,  $\mathbf{x}_1, \mathbf{x}_2 \in C$  such that for all  $\mathbf{x} \in C$ ,*

$$f(\mathbf{x}_1) \leq f(\mathbf{x}) \leq f(\mathbf{x}_2).$$

**Proof:** Let  $M = \sup \{f(\mathbf{x}) : \mathbf{x} \in C\}$ . Recall this means  $+\infty$  if  $f$  is not bounded above and it equals the least upper bound of these values of  $f$  if  $f$  is bounded above. Then there exists a sequence,  $\{\mathbf{x}_n\}$  such that  $f(\mathbf{x}_n) \rightarrow M$ . Since  $C$  is sequentially compact, there exists a subsequence,  $\mathbf{x}_{n_k}$ , and a point,  $\mathbf{x} \in C$  such that  $\mathbf{x}_{n_k} \rightarrow \mathbf{x}$ . But then since  $f$  is continuous at  $\mathbf{x}$ , it follows from Theorem 14.6.10 on Page 342 that  $f(\mathbf{x}) = \lim_{k \rightarrow \infty} f(\mathbf{x}_{n_k}) = M$ . This proves  $f$  achieves its maximum and also shows its maximum is less than  $\infty$ . Let  $\mathbf{x}_2 = \mathbf{x}$ . The case of a minimum is handled similarly.

Recall that a function is uniformly continuous if the following definition holds.

**Definition 14.9.7** *Let  $\mathbf{f} : D(\mathbf{f}) \rightarrow \mathbb{R}^q$ . Then  $\mathbf{f}$  is uniformly continuous if for every  $\varepsilon > 0$  there exists  $\delta > 0$  such that whenever  $|\mathbf{x} - \mathbf{y}| < \delta$ , it follows  $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| < \varepsilon$ .*

**Theorem 14.9.8** *Let  $\mathbf{f} : C \rightarrow \mathbb{R}^q$  be continuous where  $C$  is a closed and bounded set in  $\mathbb{R}^p$ . Then  $\mathbf{f}$  is uniformly continuous on  $C$ .*

**Proof:** If this is not so, there exists  $\varepsilon > 0$  and pairs of points,  $\mathbf{x}_n$  and  $\mathbf{y}_n$  satisfying  $|\mathbf{x}_n - \mathbf{y}_n| < 1/n$  but  $|\mathbf{f}(\mathbf{x}_n) - \mathbf{f}(\mathbf{y}_n)| \geq \varepsilon$ . Since  $C$  is sequentially compact, there exists  $\mathbf{x} \in C$  and a subsequence,  $\{\mathbf{x}_{n_k}\}$  satisfying  $\mathbf{x}_{n_k} \rightarrow \mathbf{x}$ . But  $|\mathbf{x}_{n_k} - \mathbf{y}_{n_k}| < 1/k$  and so  $\mathbf{y}_{n_k} \rightarrow \mathbf{x}$  also. Therefore, from Theorem 14.6.10 on Page 342,

$$\varepsilon \leq \lim_{k \rightarrow \infty} |\mathbf{f}(\mathbf{x}_{n_k}) - \mathbf{f}(\mathbf{y}_{n_k})| = |\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x})| = 0,$$

a contradiction. This proves the theorem.



# Limits And Derivatives

## 15.1 Limits Of A Vector Valued Function

The above discussion considered expressions like

$$\frac{\mathbf{f}(t_0 + h) - \mathbf{f}(t_0)}{h}$$

and determined what they get close to as  $h$  gets small. In other words it is desired to consider

$$\lim_{h \rightarrow 0} \frac{\mathbf{f}(t_0 + h) - \mathbf{f}(t_0)}{h}$$

Specializing to functions of one variable, one can give a meaning to

$$\lim_{s \rightarrow t+} \mathbf{f}(s), \lim_{s \rightarrow t-} \mathbf{f}(s), \lim_{s \rightarrow \infty} \mathbf{f}(s),$$

and

$$\lim_{s \rightarrow -\infty} \mathbf{f}(s).$$

**Definition 15.1.1** *In the case where  $D(\mathbf{f})$  is only assumed to satisfy  $D(\mathbf{f}) \supseteq (t, t + r)$ ,*

$$\lim_{s \rightarrow t+} \mathbf{f}(s) = \mathbf{L}$$

*if and only if for all  $\varepsilon > 0$  there exists  $\delta > 0$  such that if*

$$0 < s - t < \delta,$$

*then*

$$|\mathbf{f}(s) - \mathbf{L}| < \varepsilon.$$

*In the case where  $D(\mathbf{f})$  is only assumed to satisfy  $D(\mathbf{f}) \supseteq (t - r, t)$ ,*

$$\lim_{s \rightarrow t-} \mathbf{f}(s) = \mathbf{L}$$

*if and only if for all  $\varepsilon > 0$  there exists  $\delta > 0$  such that if*

$$0 < t - s < \delta,$$

*then*

$$|\mathbf{f}(s) - \mathbf{L}| < \varepsilon.$$

One can also consider limits as a variable “approaches” infinity. Of course nothing is “close” to infinity and so this requires a slightly different definition.

$$\lim_{t \rightarrow \infty} \mathbf{f}(t) = \mathbf{L}$$

if for every  $\varepsilon > 0$  there exists  $l$  such that whenever  $t > l$ ,

$$|\mathbf{f}(t) - \mathbf{L}| < \varepsilon \quad (15.1)$$

and

$$\lim_{t \rightarrow -\infty} \mathbf{f}(t) = \mathbf{L}$$

if for every  $\varepsilon > 0$  there exists  $l$  such that whenever  $t < l$ , (15.1) holds.

Note that in all of this the definitions are identical to the case of scalar valued functions. The only difference is that here  $|\cdot|$  refers to the norm or length in  $\mathbb{R}^p$  where maybe  $p > 1$ .

**Example 15.1.2** Let  $\mathbf{f}(t) = (\cos t, \sin t, t^2 + 1, \ln(t))$ . Find  $\lim_{t \rightarrow \pi/2} \mathbf{f}(t)$ .

Use Theorem 14.4.7 on Page 338 and the continuity of the functions to write this limit equals

$$\begin{aligned} & \left( \lim_{t \rightarrow \pi/2} \cos t, \lim_{t \rightarrow \pi/2} \sin t, \lim_{t \rightarrow \pi/2} (t^2 + 1), \lim_{t \rightarrow \pi/2} \ln(t) \right) \\ &= \left( 0, 1, \ln\left(\frac{\pi^2}{4} + 1\right), \ln\left(\frac{\pi}{2}\right) \right). \end{aligned}$$

**Example 15.1.3** Let  $\mathbf{f}(t) = \left(\frac{\sin t}{t}, t^2, t + 1\right)$ . Find  $\lim_{t \rightarrow 0} \mathbf{f}(t)$ .

Recall that  $\lim_{t \rightarrow 0} \frac{\sin t}{t} = 1$ . Then from Theorem 14.4.7 on Page 338,  $\lim_{t \rightarrow 0} \mathbf{f}(t) = (1, 0, 1)$ .

## 15.2 The Derivative And Integral

The following definition is on the derivative and integral of a vector valued function of one variable.

**Definition 15.2.1** The derivative of a function,  $\mathbf{f}'(t)$ , is defined as the following limit whenever the limit exists. If the limit does not exist, then neither does  $\mathbf{f}'(t)$ .

$$\lim_{h \rightarrow 0} \frac{\mathbf{f}(t+h) - \mathbf{f}(t)}{h} \equiv \mathbf{f}'(t)$$

The function of  $h$  on the left is called the difference quotient just as it was for a scalar valued function. If  $\mathbf{f}(t) = (f_1(t), \dots, f_p(t))$  and  $\int_a^b f_i(t) dt$  exists for each  $i = 1, \dots, p$ , then  $\int_a^b \mathbf{f}(t) dt$  is defined as the vector,

$$\left( \int_a^b f_1(t) dt, \dots, \int_a^b f_p(t) dt \right).$$

This is what is meant by saying  $\mathbf{f} \in R([a, b])$ .



This is exactly like the definition for a scalar valued function. As before,

$$\mathbf{f}'(x) = \lim_{y \rightarrow x} \frac{\mathbf{f}(y) - \mathbf{f}(x)}{y - x}.$$

As in the case of a scalar valued function, differentiability implies continuity but not the other way around.

**Theorem 15.2.2** *If  $\mathbf{f}'(t)$  exists, then  $\mathbf{f}$  is continuous at  $t$ .*

**Proof:** Suppose  $\varepsilon > 0$  is given and choose  $\delta_1 > 0$  such that if  $|h| < \delta_1$ ,

$$\left| \frac{\mathbf{f}(t+h) - \mathbf{f}(t)}{h} - \mathbf{f}'(t) \right| < 1.$$

then for such  $h$ , the triangle inequality implies

$$|\mathbf{f}(t+h) - \mathbf{f}(t)| < |h| + |\mathbf{f}'(t)| |h|.$$

Now letting  $\delta < \min\left(\delta_1, \frac{\varepsilon}{1+|\mathbf{f}'(t)|}\right)$  it follows if  $|h| < \delta$ , then

$$|\mathbf{f}(t+h) - \mathbf{f}(t)| < \varepsilon.$$

Letting  $y = h + t$ , this shows that if  $|y - t| < \delta$ ,

$$|\mathbf{f}(y) - \mathbf{f}(t)| < \varepsilon$$

which proves  $\mathbf{f}$  is continuous at  $t$ . This proves the theorem.

As in the scalar case, there is a fundamental theorem of calculus.

**Theorem 15.2.3** *If  $\mathbf{f} \in R([a, b])$  and if  $\mathbf{f}$  is continuous at  $t \in (a, b)$ , then*

$$\frac{d}{dt} \left( \int_a^t \mathbf{f}(s) \, ds \right) = \mathbf{f}(t).$$

**Proof:** Say  $\mathbf{f}(t) = (f_1(t), \dots, f_p(t))$ . Then it follows

$$\frac{1}{h} \int_a^{t+h} \mathbf{f}(s) \, ds - \frac{1}{h} \int_a^t \mathbf{f}(s) \, ds = \left( \frac{1}{h} \int_t^{t+h} f_1(s) \, ds, \dots, \frac{1}{h} \int_t^{t+h} f_p(s) \, ds \right)$$

and  $\lim_{h \rightarrow 0} \frac{1}{h} \int_t^{t+h} f_i(s) \, ds = f_i(t)$  for each  $i = 1, \dots, p$  from the fundamental theorem of calculus for scalar valued functions. Therefore,

$$\lim_{h \rightarrow 0} \frac{1}{h} \int_a^{t+h} \mathbf{f}(s) \, ds - \frac{1}{h} \int_a^t \mathbf{f}(s) \, ds = (f_1(t), \dots, f_p(t)) = \mathbf{f}(t)$$

and this proves the claim.

**Example 15.2.4** *Let  $\mathbf{f}(x) = \mathbf{c}$  where  $\mathbf{c}$  is a constant. Find  $\mathbf{f}'(x)$ .*

The difference quotient,

$$\frac{\mathbf{f}(x+h) - \mathbf{f}(x)}{h} = \frac{\mathbf{c} - \mathbf{c}}{h} = \mathbf{0}$$

Therefore,

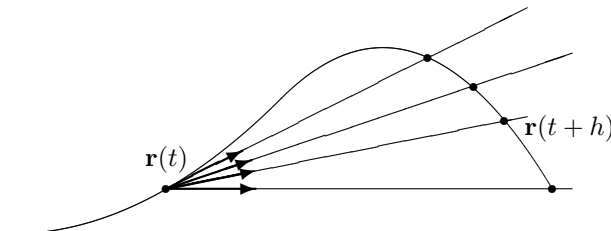
$$\lim_{h \rightarrow 0} \frac{\mathbf{f}(x+h) - \mathbf{f}(x)}{h} = \lim_{h \rightarrow 0} \mathbf{0} = \mathbf{0}$$

**Example 15.2.5** *Let  $\mathbf{f}(t) = (at, bt)$  where  $a, b$  are constants. Find  $\mathbf{f}'(t)$ .*

From the above discussion this derivative is just the vector valued functions whose components consist of the derivatives of the components of  $\mathbf{f}$ . Thus  $\mathbf{f}'(t) = (a, b)$ .

### 15.2.1 Geometric And Physical Significance Of The Derivative

Suppose  $\mathbf{r}$  is a vector valued function of a parameter,  $t$  not necessarily time and consider the following picture of the points traced out by  $\mathbf{r}$ .



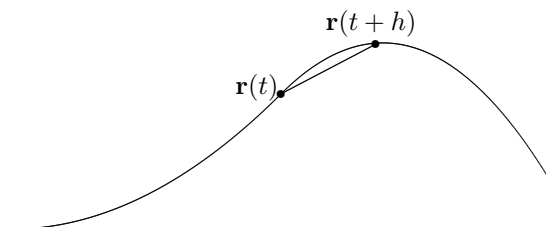
In this picture there are unit vectors in the direction of the vector from  $\mathbf{r}(t)$  to  $\mathbf{r}(t+h)$ . You can see that it is reasonable to suppose these unit vectors, if they converge, converge to a unit vector,  $\mathbf{T}$  which is tangent to the curve at the point  $\mathbf{r}(t)$ . Now each of these unit vectors is of the form

$$\frac{\mathbf{r}(t+h) - \mathbf{r}(t)}{|\mathbf{r}(t+h) - \mathbf{r}(t)|} \equiv \mathbf{T}_h.$$

Thus  $\mathbf{T}_h \rightarrow \mathbf{T}$ , a unit tangent vector to the curve at the point  $\mathbf{r}(t)$ . Therefore,

$$\begin{aligned} \mathbf{r}'(t) &\equiv \lim_{h \rightarrow 0} \frac{\mathbf{r}(t+h) - \mathbf{r}(t)}{h} = \lim_{h \rightarrow 0} \frac{|\mathbf{r}(t+h) - \mathbf{r}(t)|}{h} \frac{\mathbf{r}(t+h) - \mathbf{r}(t)}{|\mathbf{r}(t+h) - \mathbf{r}(t)|} \\ &= \lim_{h \rightarrow 0} \frac{|\mathbf{r}(t+h) - \mathbf{r}(t)|}{h} \mathbf{T}_h = |\mathbf{r}'(t)| \mathbf{T}. \end{aligned}$$

In the case that  $t$  is time, the expression  $|\mathbf{r}(t+h) - \mathbf{r}(t)|$  is a good approximation for the distance traveled by the object on the time interval  $[t, t+h]$ . The real distance would be the length of the curve joining the two points but if  $h$  is very small, this is essentially equal to  $|\mathbf{r}(t+h) - \mathbf{r}(t)|$  as suggested by the picture below.



Therefore,

$$\frac{|\mathbf{r}(t+h) - \mathbf{r}(t)|}{h}$$

gives for small  $h$ , the approximate distance travelled on the time interval,  $[t, t+h]$  divided by the length of time,  $h$ . Therefore, this expression is really the average speed of the object on this small time interval and so the limit as  $h \rightarrow 0$ , deserves to be called the instantaneous speed of the object. Thus  $|\mathbf{r}'(t)| \mathbf{T}$  represents the speed times a unit direction vector,  $\mathbf{T}$  which defines the direction in which the object is moving. Thus  $\mathbf{r}'(t)$  is the velocity of the object. This is the physical significance of the derivative when  $t$  is time.

How do you go about computing  $\mathbf{r}'(t)$ ? Letting  $\mathbf{r}(t) = (r_1(t), \dots, r_q(t))$ , the expression

$$\frac{\mathbf{r}(t_0+h) - \mathbf{r}(t_0)}{h} \tag{15.2}$$

is equal to

$$\left( \frac{r_1(t_0 + h) - r_1(t_0)}{h}, \dots, \frac{r_q(t_0 + h) - r_q(t_0)}{h} \right).$$

Then as  $h$  converges to 0, (15.2) converges to

$$\mathbf{v} \equiv (v_1, \dots, v_q)$$

where  $v_k = r'_k(t)$ . This by Theorem 14.4.7 on Page 338, which says that the term in (15.2) gets close to a vector,  $\mathbf{v}$  if and only if all the coordinate functions of the term in (15.2) get close to the corresponding coordinate functions of  $\mathbf{v}$ .

In the case where  $t$  is time, this simply says the velocity vector equals the vector whose components are the derivatives of the components of the displacement vector,  $\mathbf{r}(t)$ .

In any case, the vector,  $\mathbf{T}$  determines a direction vector which is tangent to the curve at the point,  $\mathbf{r}(t)$  and so it is possible to find parametric equations for the line tangent to the curve at various points.

**Example 15.2.6** Let  $\mathbf{r}(t) = (\sin t, t^2, t + 1)$  for  $t \in [0, 5]$ . Find a tangent line to the curve parameterized by  $\mathbf{r}$  at the point  $\mathbf{r}(2)$ .

From the above discussion, a direction vector has the same direction as  $\mathbf{r}'(2)$ . Therefore, it suffices to simply use  $\mathbf{r}'(2)$  as a direction vector for the line.  $\mathbf{r}'(2) = (\cos 2, 4, 1)$ . Therefore, a parametric equation for the tangent line is

$$(\sin 2, 4, 3) + t(\cos 2, 4, 1) = (x, y, z).$$

**Example 15.2.7** Let  $\mathbf{r}(t) = (\sin t, t^2, t + 1)$  for  $t \in [0, 5]$ . Find the velocity vector when  $t = 1$ .

From the above discussion, this is simply  $\mathbf{r}'(1) = (\cos 1, 2, 1)$ .

## 15.2.2 Differentiation Rules

There are rules which relate the derivative to the various operations done with vectors such as the dot product, the cross product, and vector addition and scalar multiplication.

**Theorem 15.2.8** Let  $a, b \in \mathbb{R}$  and suppose  $\mathbf{f}'(t)$  and  $\mathbf{g}'(t)$  exist. Then the following formulas are obtained.

$$(a\mathbf{f} + b\mathbf{g})'(t) = a\mathbf{f}'(t) + b\mathbf{g}'(t). \quad (15.3)$$

$$(\mathbf{f} \cdot \mathbf{g})'(t) = \mathbf{f}'(t) \cdot \mathbf{g}(t) + \mathbf{f}(t) \cdot \mathbf{g}'(t) \quad (15.4)$$

If  $\mathbf{f}, \mathbf{g}$  have values in  $\mathbb{R}^3$ , then

$$(\mathbf{f} \times \mathbf{g})'(t) = \mathbf{f}(t) \times \mathbf{g}'(t) + \mathbf{f}'(t) \times \mathbf{g}(t) \quad (15.5)$$

The formulas, (15.4), and (15.5) are referred to as the product rule.

**Proof:** The first formula is left for you to prove. Consider the second, (15.4).

$$\begin{aligned}
 \lim_{h \rightarrow 0} \frac{\mathbf{f} \cdot \mathbf{g}(t+h) - \mathbf{f}\mathbf{g}(t)}{h} &= \lim_{h \rightarrow 0} \frac{\mathbf{f}(t+h) \cdot \mathbf{g}(t+h) - \mathbf{f}(t+h) \cdot \mathbf{g}(t)}{h} + \frac{\mathbf{f}(t+h) \cdot \mathbf{g}(t) - \mathbf{f}(t) \cdot \mathbf{g}(t)}{h} \\
 &= \lim_{h \rightarrow 0} \left( \mathbf{f}(t+h) \cdot \frac{(\mathbf{g}(t+h) - \mathbf{g}(t))}{h} + \frac{(\mathbf{f}(t+h) - \mathbf{f}(t))}{h} \cdot \mathbf{g}(t) \right) \\
 &= \lim_{h \rightarrow 0} \sum_{k=1}^n f_k(t+h) \frac{(g_k(t+h) - g_k(t))}{h} + \sum_{k=1}^n \frac{(f_k(t+h) - f_k(t))}{h} g_k(t) \\
 &= \sum_{k=1}^n f_k(t) g'_k(t) + \sum_{k=1}^n f'_k(t) g_k(t) \\
 &= \mathbf{f}'(t) \cdot \mathbf{g}(t) + \mathbf{f}(t) \cdot \mathbf{g}'(t).
 \end{aligned}$$

Formula (15.5) is left as an exercise which follows from the product rule and the definition of the cross product in terms of components given on Page 323.

**Example 15.2.9** Let

$$\mathbf{r}(t) = (t^2, \sin t, \cos t)$$

and let  $\mathbf{p}(t) = (t, \ln(t+1), 2t)$ . Find  $(\mathbf{r}(t) \times \mathbf{p}(t))'$ .

From (15.5) this equals  $(2t, \cos t, -\sin t) \times (t, \ln(t+1), 2t) + (t^2, \sin t, \cos t) \times \left(1, \frac{1}{t+1}, 2\right)$ .

**Example 15.2.10** Let  $\mathbf{r}(t) = (t^2, \sin t, \cos t)$  Find  $\int_0^\pi \mathbf{r}(t) dt$ .

This equals  $(\int_0^\pi t^2 dt, \int_0^\pi \sin t dt, \int_0^\pi \cos t dt) = (\frac{1}{3}\pi^3, 2, 0)$ .

**Example 15.2.11** An object has position  $\mathbf{r}(t) = \left(t^3, \frac{t}{1+t}, \sqrt{t^2+2}\right)$  kilometers where  $t$  is given in hours. Find the velocity of the object in kilometers per hour when  $t = 1$ .

Recall the velocity at time  $t$  was  $\mathbf{r}'(t)$ . Therefore, find  $\mathbf{r}'(t)$  and plug in  $t = 1$  to find the velocity.

$$\begin{aligned}
 \mathbf{r}'(t) &= \left( 3t^2, \frac{1(1+t) - t}{(1+t)^2}, \frac{1}{2}(t^2+2)^{-1/2} 2t \right) \\
 &= \left( 3t^2, \frac{1}{(1+t)^2}, \frac{1}{\sqrt{(t^2+2)}} t \right)
 \end{aligned}$$

When  $t = 1$ , the velocity is

$$\mathbf{r}'(1) = \left( 3, \frac{1}{4}, \frac{1}{\sqrt{3}} \right) \text{ kilometers per hour.}$$

Obviously, this can be continued. That is, you can consider the possibility of taking the derivative of the derivative and then the derivative of that and so forth. The main thing to consider about this is the notation and it is exactly like it was in the case of a scalar valued function presented earlier. Thus  $\mathbf{r}''(t)$  denotes the second derivative.

When you are given a vector valued function of one variable, sometimes it is possible to give a simple description of the curve which results. Usually it is not possible to do this!

**Example 15.2.12** Describe the curve which results from the vector valued function,  $\mathbf{r}(t) = (\cos 2t, \sin 2t, t)$  where  $t \in \mathbb{R}$ .

The first two components indicate that for  $\mathbf{r}(t) = (x(t), y(t), z(t))$ , the pair,  $(x(t), y(t))$  traces out a circle. While it is doing so,  $z(t)$  is moving at a steady rate in the positive direction. Therefore, the curve which results is a cork skrew shaped thing called a helix.

As an application of the theorems for differentiating curves, here is an interesting application. It is also a situation where the curve can be identified as something familiar.

**Example 15.2.13** *Sound waves have the angle of incidence equal to the angle of reflection. Suppose you are in a large room and you make a sound. The sound waves spread out and you would expect your sound to be inaudible very far away. But what if the room were shaped so that the sound is reflected off the wall toward a single point, possibly far away from you? Then you might have the interesting phenomenon of someone far away hearing what you said quite clearly. How should the room be designed?*

Suppose you are located at the point  $\mathbf{P}_0$  and the point where your sound is to be reflected is  $\mathbf{P}_1$ . Consider a plane which contains the two points and let  $\mathbf{r}(t)$  denote a parameterization of the intersection of this plane with the walls of the room. Then the condition that the angle of reflection equals the angle of incidence reduces to saying the angle between  $\mathbf{P}_0 - \mathbf{r}(t)$  and  $-\mathbf{r}'(t)$  equals the angle between  $\mathbf{P}_1 - \mathbf{r}(t)$  and  $\mathbf{r}'(t)$ . Draw a picture to see this. Therefore,

$$\frac{(\mathbf{P}_0 - \mathbf{r}(t)) \cdot (-\mathbf{r}'(t))}{|\mathbf{P}_0 - \mathbf{r}(t)| |\mathbf{r}'(t)|} = \frac{(\mathbf{P}_1 - \mathbf{r}(t)) \cdot (\mathbf{r}'(t))}{|\mathbf{P}_1 - \mathbf{r}(t)| |\mathbf{r}'(t)|}.$$

This reduces to

$$\frac{(\mathbf{r}(t) - \mathbf{P}_0) \cdot (-\mathbf{r}'(t))}{|\mathbf{r}(t) - \mathbf{P}_0|} = \frac{(\mathbf{r}(t) - \mathbf{P}_1) \cdot (\mathbf{r}'(t))}{|\mathbf{r}(t) - \mathbf{P}_1|} \quad (15.6)$$

Now

$$\frac{(\mathbf{r}(t) - \mathbf{P}_1) \cdot (\mathbf{r}'(t))}{|\mathbf{r}(t) - \mathbf{P}_1|} = \frac{d}{dt} |\mathbf{r}(t) - \mathbf{P}_1|$$

and a similar formula holds for  $\mathbf{P}_1$  replaced with  $\mathbf{P}_0$ . This is because

$$|\mathbf{r}(t) - \mathbf{P}_1| = \sqrt{(\mathbf{r}(t) - \mathbf{P}_1) \cdot (\mathbf{r}(t) - \mathbf{P}_1)}$$

and so using the chain rule and product rule,

$$\begin{aligned} \frac{d}{dt} |\mathbf{r}(t) - \mathbf{P}_1| &= \frac{1}{2} ((\mathbf{r}(t) - \mathbf{P}_1) \cdot (\mathbf{r}(t) - \mathbf{P}_1))^{-1/2} 2 ((\mathbf{r}(t) - \mathbf{P}_1) \cdot \mathbf{r}'(t)) \\ &= \frac{(\mathbf{r}(t) - \mathbf{P}_1) \cdot (\mathbf{r}'(t))}{|\mathbf{r}(t) - \mathbf{P}_1|}. \end{aligned}$$

Therefore, from (15.6),

$$\frac{d}{dt} (|\mathbf{r}(t) - \mathbf{P}_1|) + \frac{d}{dt} (|\mathbf{r}(t) - \mathbf{P}_0|) = 0$$

showing that  $|\mathbf{r}(t) - \mathbf{P}_1| + |\mathbf{r}(t) - \mathbf{P}_0| = C$  for some constant,  $C$ . This implies the curve of intersection of the plane with the room is an ellipse having  $\mathbf{P}_0$  and  $\mathbf{P}_1$  as the foci.

## 15.3 Leibniz's Notation

Leibniz's notation also generalizes routinely. For example,  $\frac{dy}{dt} = \mathbf{y}'(t)$  with other similar notations holding.

## 15.4 Exercises

1. Find the following limits if possible

(a)  $\lim_{x \rightarrow 0+} \left( \frac{|x|}{x}, \sin x/x, \cos x \right)$

(b)  $\lim_{x \rightarrow 0+} \left( \frac{x}{|x|}, \sec x, e^x \right)$

(c)  $\lim_{x \rightarrow 4} \left( \frac{x^2-16}{x+4}, x+7, \frac{\tan 4x}{5x} \right)$

(d)  $\lim_{x \rightarrow \infty} \left( \frac{x}{1+x^2}, \frac{x^2}{1+x^2}, \frac{\sin x^2}{x} \right)$

2. Find  $\lim_{x \rightarrow 2} \left( \frac{x^2-4}{x+2}, x^2+2x-1, \frac{x^2-4}{x-2} \right)$ .

3. Prove from the definition that  $\lim_{x \rightarrow a} (\sqrt[3]{x}, x+1) = (\sqrt[3]{a}, a+1)$  for all  $a \in \mathbb{R}$ . **Hint:** You might want to use the formula for the difference of two cubes,

$$a^3 - b^3 = (a - b)(a^2 + ab + b^2).$$

4. Let  $\mathbf{r}(t) = \left( 4 + (t-1)^2, \sqrt{t^2+1}(t-1)^3, \frac{(t-1)^3}{t^5} \right)$  describe the position of an object in  $\mathbb{R}^3$  as a function of  $t$  where  $t$  is measured in seconds and  $\mathbf{r}(t)$  is measured in meters. Is the velocity of this object ever equal to zero? If so, find the value of  $t$  at which this occurs and the point in  $\mathbb{R}^3$  at which the velocity is zero.

5. Let  $\mathbf{r}(t) = (\sin 2t, t^2, 2t+1)$  for  $t \in [0, 4]$ . Find a tangent line to the curve parameterized by  $\mathbf{r}$  at the point  $\mathbf{r}(2)$ .

6. Let  $\mathbf{r}(t) = (t, \sin t^2, t+1)$  for  $t \in [0, 5]$ . Find a tangent line to the curve parameterized by  $\mathbf{r}$  at the point  $\mathbf{r}(2)$ .

7. Let  $\mathbf{r}(t) = (\sin t, t^2, \cos(t^2))$  for  $t \in [0, 5]$ . Find a tangent line to the curve parameterized by  $\mathbf{r}$  at the point  $\mathbf{r}(2)$ .

8. Let  $\mathbf{r}(t) = (\sin t, \cos(t^2), t+1)$  for  $t \in [0, 5]$ . Find the velocity when  $t = 3$ .

9. Let  $\mathbf{r}(t) = (\sin t, t^2, t+1)$  for  $t \in [0, 5]$ . Find the velocity when  $t = 3$ .

10. Let  $\mathbf{r}(t) = (t, \ln(t^2+1), t+1)$  for  $t \in [0, 5]$ . Find the velocity when  $t = 3$ .

11. Suppose an object has position  $\mathbf{r}(t) \in \mathbb{R}^3$  where  $\mathbf{r}$  is differentiable and suppose also that  $|\mathbf{r}(t)| = c$  where  $c$  is a constant.

- (a) Show first that this condition does not require  $\mathbf{r}(t)$  to be a constant. **Hint:** You can do this either mathematically or by giving a physical example.

- (b) Show that you can conclude that  $\mathbf{r}'(t) \cdot \mathbf{r}(t) = 0$ . That is, the velocity is always perpendicular to the displacement.

12. Prove (15.5) from the component description of the cross product.

13. Prove (15.5) from the formula  $(\mathbf{f} \times \mathbf{g})_i = \varepsilon_{ijk} f_j g_k$ .

14. Prove (15.5) directly from the definition of the derivative without considering components.

15. A bezier curve in  $\mathbb{R}^n$  is a vector valued function of the form

$$\mathbf{y}(t) = \sum_{k=0}^n \binom{n}{k} \mathbf{x}_k (1-t)^{n-k} t^k$$

where here the  $\binom{n}{k}$  are the binomial coefficients and  $\mathbf{x}_k$  are  $n+1$  points in  $\mathbb{R}^n$ . Show that  $\mathbf{y}(0) = \mathbf{x}_0$ ,  $\mathbf{y}(1) = \mathbf{x}_n$ , and find  $\mathbf{y}'(0)$  and  $\mathbf{y}'(1)$ . Recall that  $\binom{n}{0} = \binom{n}{n} = 1$  and  $\binom{n}{n-1} = \binom{n}{1} = n$ . Curves of this sort are important in various computer programs.

16. Suppose  $\mathbf{r}(t)$ ,  $\mathbf{s}(t)$ , and  $\mathbf{p}(t)$  are three differentiable functions of  $t$  which have values in  $\mathbb{R}^3$ . Find a formula for  $(\mathbf{r}(t) \times \mathbf{s}(t) \cdot \mathbf{p}(t))'$ .
17. If  $\mathbf{r}'(t) = \mathbf{0}$  for all  $t \in (a, b)$ , show there exists a constant vector,  $\mathbf{c}$  such that  $\mathbf{r}(t) = \mathbf{c}$  for all  $t \in (a, b)$ .
18. If  $\mathbf{F}'(t) = \mathbf{f}(t)$  for all  $t \in (a, b)$  and  $\mathbf{F}$  is continuous on  $[a, b]$ , show  $\int_a^b \mathbf{f}(t) dt = \mathbf{F}(b) - \mathbf{F}(a)$ .

## 15.5 Newton's Laws Of Motion

**Definition 15.5.1** Let  $\mathbf{r}(t)$  denote the position of an object. Then the acceleration of the object is defined to be  $\mathbf{r}''(t)$ .

Newton's<sup>1</sup> first law is: "Every body persists in its state of rest or of uniform motion in a straight line unless it is compelled to change that state by forces impressed on it."

Newton's second law is:

$$\mathbf{F} = m\mathbf{a} \tag{15.7}$$

where  $\mathbf{a}$  is the acceleration and  $m$  is the mass of the object.

Newton's third law states: "To every action there is always opposed an equal reaction; or, the mutual actions of two bodies upon each other are always equal, and directed to contrary parts."

Of these laws, only the second two are independent of each other, the first law being implied by the second. The third law says roughly that if you apply a force to something, the thing applies the same force back.

The second law is the one of most interest. Note that the statement of this law depends on the concept of the derivative because the acceleration is defined as a derivative. Newton used calculus and these laws to solve profound problems involving the motion of the planets and other problems in mechanics. The next example involves the concept that if you know the force along with the initial velocity and initial position, then you can determine the position.

**Example 15.5.2** Let  $\mathbf{r}(t)$  denote the position of an object of mass 2 kilogram at time  $t$  and suppose the force acting on the object is given by  $\mathbf{F}(t) = (t, 1-t^2, 2e^{-t})$ . Suppose  $\mathbf{r}(0) = (1, 0, 1)$  meters, and  $\mathbf{r}'(0) = (0, 1, 1)$  meters/sec. Find  $\mathbf{r}(t)$ .

---

<sup>1</sup>Isaac Newton 1642-1727 is often credited with inventing calculus although this is not correct since most of the ideas were in existence earlier. However, he made major contributions to the subject partly in order to study physics and astronomy. He formulated the laws of gravity, made major contributions to optics, and stated the fundamental laws of mechanics listed here. He invented a version of the binomial theorem when he was only 23 years old and built a reflecting telescope. He showed that Kepler's laws for the motion of the planets came from calculus and his laws of gravitation. In 1686 he published an important book, Principia, in which many of his ideas are found. Newton was also very interested in theology and had strong views on the nature of God which were based on his study of the Bible and early Christian writings. He finished his life as Master of the Mint.

By Newton's second law,  $2\mathbf{r}''(t) = \mathbf{F}(t) = (t, 1 - t^2, 2e^{-t})$  and so

$$\mathbf{r}''(t) = (t/2, (1 - t^2)/2, e^{-t}).$$

Therefore the velocity is given by

$$\mathbf{r}'(t) = \left( \frac{t^2}{4}, \frac{t - t^3/3}{2}, -e^{-t} \right) + \mathbf{c}$$

where  $\mathbf{c}$  is a constant vector which must be determined from the initial condition given for the velocity. Thus letting  $\mathbf{c} = (c_1, c_2, c_3)$ ,

$$(0, 1, 1) = (0, 0, -1) + (c_1, c_2, c_3)$$

which requires  $c_1 = 0$ ,  $c_2 = 1$ , and  $c_3 = 2$ . Therefore, the velocity is found.

$$\mathbf{r}'(t) = \left( \frac{t^2}{4}, \frac{t - t^3/3}{2} + 1, -e^{-t} + 2 \right).$$

Now from this, the displacement must equal

$$\mathbf{r}(t) = \left( \frac{t^3}{12}, \frac{t^2/2 - t^4/12}{2} + t, e^{-t} + 2t \right) + (C_1, C_2, C_3)$$

where the constant vector,  $(C_1, C_2, C_3)$  must be determined from the initial condition for the displacement. Thus

$$\mathbf{r}(0) = (1, 0, 1) = (0, 0, 1) + (C_1, C_2, C_3)$$

which means  $C_1 = 1$ ,  $C_2 = 0$ , and  $C_3 = 0$ . Therefore, the displacement has also been found.

$$\mathbf{r}(t) = \left( \frac{t^3}{12} + 1, \frac{t^2/2 - t^4/12}{2} + t, e^{-t} + 2t \right) \text{ meters.}$$

Actually, in applications of this sort of thing acceleration does not usually come to you as a nice given function written in terms of simple functions you understand. Rather, it comes as measurements taken by instruments and the position is continuously being updated based on this information. Another situation which often occurs is the case when the forces on the object depend not just on time but also on the position or velocity of the object.

**Example 15.5.3** *An artillery piece is fired at ground level on a level plain. The angle of elevation is  $\pi/6$  radians and the speed of the shell is 400 meters per second. How far does the shell fly before hitting the ground?*

Neglect air resistance in this problem. Also let the direction of flight be along the positive  $x$  axis. Thus the initial velocity is the vector,  $400 \cos(\pi/6) \mathbf{i} + 400 \sin(\pi/6) \mathbf{j}$  while the only force experienced by the shell after leaving the artillery piece is the force of gravity,  $-mg\mathbf{j}$  where  $m$  is the mass of the shell. The acceleration of gravity equals 9.8 meters per sec<sup>2</sup> and so the following needs to be solved.

$$m\mathbf{r}''(t) = -mg\mathbf{j}, \mathbf{r}(0) = (0, 0), \mathbf{r}'(0) = 400 \cos(\pi/6) \mathbf{i} + 400 \sin(\pi/6) \mathbf{j}.$$

Denoting  $\mathbf{r}(t)$  as  $(x(t), y(t))$ ,

$$x''(t) = 0, y''(t) = -g.$$



Therefore,  $y'(t) = -gt + C$  and from the information on the initial velocity,  $C = 400 \sin(\pi/6) = 200$ . Thus

$$y(t) = -4.9t^2 + 200t + D.$$

$D = 0$  because the artillery piece is fired at ground level which requires both  $x$  and  $y$  to equal zero at this time. Similarly,  $x'(t) = 400 \cos(\pi/6)$  so  $x(t) = 400 \cos(\pi/6)t = 200\sqrt{3}t$ . The shell hits the ground when  $y = 0$  and this occurs when  $-4.9t^2 + 200t = 0$ . Thus  $t = 40.8163265306$  seconds and so at this time,

$$x = 200\sqrt{3}(40.8163265306) = 14139.1902659 \text{ meters.}$$

The next example is more complicated because it also takes in to account air resistance.

**Example 15.5.4** *A lump of "blue ice" escapes the lavatory of a jet flying at 600 miles per hour at an altitude of 30,000 feet. This blue ice weighs 64 pounds near the earth and experiences a force of air resistance equal to  $(-.1)\mathbf{r}'(t)$  pounds. Find the position and velocity of the blue ice as a function of time measured in seconds. Also find the velocity when the lump hits the ground. Such lumps have been known to surprise people on the ground.*

The first thing needed is to obtain information which involves consistent units. The blue ice weighs 32 pounds near the earth. Thus 32 pounds is the force exerted by gravity on the lump and so its mass must be given by Newton's second law as follows.

$$64 = m \times 32.$$

Thus  $m = 2$  slugs. The slug is the unit of mass in the system involving feet and pounds. The jet is flying at 600 miles per hour. I want to change this to feet per second. Thus it flies at

$$\frac{600 \times 5280}{60 \times 60} = 880 \text{ feet per second.}$$

The explanation for this is that there are 5280 feet in a mile and so it goes  $600 \times 5280$  feet in one hour. There are  $60 \times 60$  seconds in an hour. The position of the lump of blue ice will be computed from a point on the ground directly beneath the airplane at the instant the blue ice escapes and regard the airplane as moving in the direction of the positive  $x$  axis. Thus the initial displacement is

$$\mathbf{r}(0) = (0, 30000) \text{ feet}$$

and the initial velocity is

$$\mathbf{r}'(0) = (880, 0) \text{ feet/sec.}$$

The force of gravity is

$$(0, -64) \text{ pounds}$$

and the force due to air resistance is

$$(-.1)\mathbf{r}'(t) \text{ pounds.}$$

Newton's second law yields the following initial value problem for  $\mathbf{r}(t) = (r_1(t), r_2(t))$ .

$$\begin{aligned} 2(r_1''(t), r_2''(t)) &= (-.1)(r_1'(t), r_2'(t)) + (0, -64), \quad (r_1(0), r_2(0)) = (0, 30000), \\ (r_1'(0), r_2'(0)) &= (880, 0) \end{aligned}$$

Therefore,

$$\begin{aligned} 2r_1''(t) + (.1)r_1'(t) &= 0 \\ 2r_2''(t) + (.1)r_2'(t) &= -64 \\ r_1(0) &= 0 \\ r_2(0) &= 30000 \\ r_1'(0) &= 880 \\ r_2'(0) &= 0 \end{aligned} \quad (15.8)$$

To save on repetition solve

$$mr'' + kr' = c, r(0) = u, r'(0) = v.$$

Divide the differential equation by  $m$  and get

$$r'' + (k/m)r' = c/m.$$

Now multiply both sides by  $e^{(k/m)t}$ . You should check this gives

$$\frac{d}{dt} \left( e^{(k/m)t} r' \right) = (c/m) e^{(k/m)t}$$

Therefore,

$$e^{(k/m)t} r' = \frac{1}{k} e^{\frac{k}{m}t} c + C$$

and using the initial condition,  $v = c/k + C$  and so

$$r'(t) = (c/k) + (v - (c/k)) e^{-\frac{k}{m}t}$$

Now this implies

$$r(t) = (c/k)t - \frac{1}{k} m e^{-\frac{k}{m}t} \left( v - \frac{c}{k} \right) + D \quad (15.9)$$

where  $D$  is a constant to be determined from the initial conditions. Thus

$$u = -\frac{m}{k} \left( v - \frac{c}{k} \right) + D$$

and so

$$r(t) = (c/k)t - \frac{1}{k} m e^{-\frac{k}{m}t} \left( v - \frac{c}{k} \right) + \left( u + \frac{m}{k} \left( v - \frac{c}{k} \right) \right).$$

Now apply this to the system (15.8) to find

$$\begin{aligned} r_1(t) &= -\frac{1}{(.1)} 2 \left( \exp \left( -\frac{(.1)}{2} t \right) \right) (880) + \left( \frac{2}{(.1)} (880) \right) \\ &= -17600.0 \exp(-.05t) + 17600.0 \end{aligned}$$

and

$$\begin{aligned} r_2(t) &= (-64/ (.1)) t - \frac{1}{(.1)} 2 \left( \exp \left( -\frac{(.1)}{2} t \right) \right) \left( \frac{64}{(.1)} \right) + \left( 30000 + \frac{2}{(.1)} \left( \frac{64}{(.1)} \right) \right) \\ &= -640.0t - 12800.0 \exp(-.05t) + 42800.0 \end{aligned}$$

This gives the coordinates of the position. What of the velocity? Using (15.9) in the same way to obtain the velocity,

$$\begin{aligned} r_1'(t) &= 880.0 \exp(-.05t), \\ r_2'(t) &= -640.0 + 640.0 \exp(-.05t). \end{aligned} \quad (15.10)$$

To determine the velocity when the blue ice hits the ground, it is necessary to find the value of  $t$  when this event takes place and then to use (15.10) to determine the velocity. It hits ground when  $r_2(t) = 0$ . Thus it suffices to solve the equation,

$$0 = -640.0t - 12800.0 \exp(-.05t) + 42800.0.$$

This is a fairly hard equation to solve using the methods of algebra. In fact, I do not have a good way to find this value of  $t$  using algebra. However if plugging in various values of  $t$  using a calculator you eventually find that when  $t = 66.14$ ,

$$-640.0(66.14) - 12800.0 \exp(-.05(66.14)) + 42800.0 = 1.588 \text{ feet.}$$

This is close enough to hitting the ground and so plugging in this value for  $t$  yields the approximate velocity,

$$(880.0 \exp(-.05(66.14)), -640.0 + 640.0 \exp(-.05(66.14))) = (32.23, -616.56).$$

Notice how because of air resistance the component of velocity in the horizontal direction is only about 32 feet per second even though this component started out at 880 feet per second while the component in the vertical direction is -616 feet per second even though this component started off at 0 feet per second. You see that air resistance can be very important so it is not enough to pretend, as is often done in beginning physics courses that everything takes place in a vacuum. Actually, this problem used several physical simplifications. It was assumed the force acting on the lump of blue ice by gravity was constant. This is not really true because it actually depends on the distance between the center of mass of the earth and the center of mass of the lump. It was also assumed the air resistance is proportional to the velocity. This is an over simplification when high speeds are involved. However, increasingly correct models can be studied in a systematic way as above.

### 15.5.1 Kinetic Energy

Newton's second law is also the basis for the notion of kinetic energy. When a force is exerted on an object which causes the object to move, it follows that the force is doing work which manifests itself in a change of velocity of the object. How is the total work done on the object by the force related to the final velocity of the object? By Newton's second law, and letting  $\mathbf{v}$  be the velocity,

$$\mathbf{F}(t) = m\mathbf{v}'(t).$$

Now in a small increment of time,  $(t, t + dt)$ , the work done on the object would be approximately equal to

$$dW = \mathbf{F}(t) \cdot \mathbf{v}(t) dt. \quad (15.11)$$

If no work has been done at time  $t = 0$ , then (15.11) implies

$$\frac{dW}{dt} = \mathbf{F} \cdot \mathbf{v}, \quad W(0) = 0.$$

Hence,

$$\frac{dW}{dt} = m\mathbf{v}'(t) \cdot \mathbf{v}(t) = \frac{m}{2} \frac{d}{dt} |\mathbf{v}(t)|^2.$$

Therefore, the total work done up to time  $t$  would be  $W(t) = \frac{m}{2} |\mathbf{v}(t)|^2 - \frac{m}{2} |\mathbf{v}_0|^2$  where  $|\mathbf{v}_0|$  denotes the initial speed of the object. This difference represents the change in the kinetic energy.

## 15.5.2 Impulse And Momentum

Work and energy involve a force acting on an object for some distance. Impulse involves a force which acts on an object for an interval of time.

**Definition 15.5.5** Let  $\mathbf{F}$  be a force which acts on an object during the time interval,  $[a, b]$ . The impulse of this force is

$$\int_a^b \mathbf{F}(t) dt.$$

This is defined as

$$\left( \int_a^b F_1(t) dt, \int_a^b F_2(t) dt, \int_a^b F_3(t) dt \right).$$

The linear momentum of an object of mass  $m$  and velocity  $\mathbf{v}$  is defined as

$$\text{Linear momentum} = m\mathbf{v}.$$

The notion of impulse and momentum are related in the following theorem.

**Theorem 15.5.6** Let  $\mathbf{F}$  be a force acting on an object of mass  $m$ . Then the impulse equals the change in momentum. More precisely,

$$\int_a^b \mathbf{F}(t) dt = m\mathbf{v}(b) - m\mathbf{v}(a).$$

**Proof:** This is really just the fundamental theorem of calculus and Newton's second law applied to the components of  $\mathbf{F}$ .

$$\int_a^b \mathbf{F}(t) dt = \int_a^b m \frac{d\mathbf{v}}{dt} dt = m\mathbf{v}(b) - m\mathbf{v}(a) \quad (15.12)$$

Now suppose two point masses,  $A$  and  $B$  collide. Newton's third law says the force exerted by mass  $A$  on mass  $B$  is equal in magnitude but opposite in direction to the force exerted by mass  $B$  on mass  $A$ . Letting the collision take place in the time interval,  $[a, b]$  and denoting the two masses by  $m_A$  and  $m_B$  and their velocities by  $\mathbf{v}_A$  and  $\mathbf{v}_B$  it follows that

$$m_A \mathbf{v}_A(b) - m_A \mathbf{v}_A(a) = \int_a^b (\text{Force of } B \text{ on } A) dt$$

and

$$\begin{aligned} m_B \mathbf{v}_B(b) - m_B \mathbf{v}_B(a) &= \int_a^b (\text{Force of } A \text{ on } B) dt \\ &= - \int_a^b (\text{Force of } B \text{ on } A) dt \\ &= - (m_A \mathbf{v}_A(b) - m_A \mathbf{v}_A(a)) \end{aligned}$$

and this shows

$$m_B \mathbf{v}_B(b) + m_A \mathbf{v}_A(b) = m_B \mathbf{v}_B(a) + m_A \mathbf{v}_A(a).$$

In other words, in a collision between two masses the total linear momentum before the collision equals the total linear momentum after the collision. This is known as the conservation of linear momentum.

## 15.6 Exercises

1. Show the solution to  $\mathbf{v}' + r\mathbf{v} = \mathbf{c}$  with the initial condition,  $\mathbf{v}(0) = \mathbf{v}_0$  is  $\mathbf{v}(t) = (\mathbf{v}_0 - \frac{\mathbf{c}}{r})e^{-rt} + (\mathbf{c}/r)$ . If  $\mathbf{v}$  is velocity and  $r = k/m$  where  $k$  is a constant for air resistance and  $m$  is the mass, and  $\mathbf{c} = \mathbf{f}/m$ , argue from Newton's second law that this is the equation for finding the velocity,  $\mathbf{v}$  of an object acted on by air resistance proportional to the velocity and a constant force,  $\mathbf{f}$ , possibly from gravity. Does there exist a terminal velocity? What is it?
2. Verify Formula (15.12) carefully by considering the components.
3. Suppose that the air resistance is proportional to the velocity but it is desired to find the constant of proportionality. Describe how you could find this constant.
4. Suppose an object having mass equal to 5 kilograms experiences a time dependent force,  $\mathbf{F}(t) = e^{-t}\mathbf{i} + \cos(t)\mathbf{j} + t^2\mathbf{k}$  meters per sec<sup>2</sup>. Suppose also that the object is at the point  $(0, 1, 1)$  meters at time  $t = 0$  and that its initial velocity at this time is  $\mathbf{v} = \mathbf{i} + \mathbf{j} - \mathbf{k}$  meters per sec. Find the position of the object as a function of  $t$ .
5. Fill in the details for the derivation of kinetic energy. In particular verify that  $m\mathbf{v}'(t) \cdot \mathbf{v}(t) = \frac{m}{2} \frac{d}{dt} |\mathbf{v}(t)|^2$ . Also, why would  $dW = \mathbf{F}(t) \cdot \mathbf{v}(t) dt$ ?
6. Suppose the force acting on an object,  $\mathbf{F}$  is always perpendicular to the velocity of the object. Thus  $\mathbf{F} \cdot \mathbf{v} = 0$ . Show the Kinetic energy of the object is constant. Such forces are sometimes called forces of constraint because they do not contribute to the speed of the object, only its direction.
7. A cannon is fired at an angle,  $\theta$  from ground level on a vast plain. The speed of the ball as it leaves the mouth of the cannon is known to be  $s$  meters per second. Neglecting air resistance, find a formula for how far the cannon ball goes before hitting the ground. Show the maximum range for the cannon ball is achieved when  $\theta = \pi/4$ .
8. Suppose in the context of Problem 7 that the cannon ball has mass 10 kilograms and it experiences a force of air resistance which is  $.01\mathbf{v}$  Newtons where  $\mathbf{v}$  is the velocity in meters per second. The acceleration of gravity is 9.8 meters per sec<sup>2</sup>. Also suppose that the initial speed is 100 meters per second. Find a formula for the displacement,  $\mathbf{r}(t)$  of the cannon ball. If the angle of elevation equals  $\pi/4$ , use a calculator or other means to estimate the time before the cannon ball hits the ground.
9. Show using Corollary 6.8.4 that Newton's first law can be obtained from the second law.
10. Show that if  $\mathbf{v}'(t) = \mathbf{0}$ , for all  $t \in (a, b)$ , then there exists a constant vector,  $\mathbf{z}$  independent of  $t$  such that  $\mathbf{v}(t) = \mathbf{z}$  for all  $t$ .
11. Suppose an object moves in three dimensional space in such a way that the only force acting on the object is directed toward a single fixed point in three dimensional space. Verify that the motion of the object takes place in a plane. **Hint:** Let  $\mathbf{r}(t)$  denote the position vector of the object from the fixed point. Then the force acting on the object must be of the form  $g(\mathbf{r}(t))\mathbf{r}(t)$  and by Newton's second law, this equals  $m\mathbf{r}''(t)$ . Therefore,

$$m\mathbf{r}'' \times \mathbf{r} = g(\mathbf{r})\mathbf{r} \times \mathbf{r} = \mathbf{0}.$$

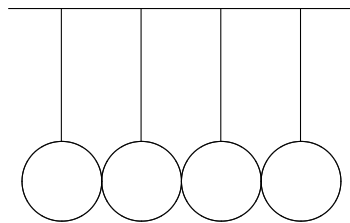
Now argue that  $\mathbf{r}'' \times \mathbf{r} = (\mathbf{r}' \times \mathbf{r})'$ , showing that  $(\mathbf{r}' \times \mathbf{r})$  must equal a constant vector,  $\mathbf{z}$ . Therefore, what can be said about  $\mathbf{z}$  and  $\mathbf{r}$ ?

12. Suppose the only forces acting on an object are the force of gravity,  $-mg\mathbf{k}$  and a force,  $\mathbf{F}$  which is perpendicular to the motion of the object. Thus  $\mathbf{F} \cdot \mathbf{v} = 0$ . Show the total energy of the object,

$$E \equiv \frac{1}{2}m|\mathbf{v}|^2 + mgz$$

is constant. Here  $\mathbf{v}$  is the velocity and the first term is the kinetic energy while the second is the potential energy. **Hint:** Use Newton's second law to show the time derivative of the above expression equals zero.

13. Using Problem 12, suppose an object slides down a frictionless inclined plane from a height of 100 feet. When it reaches the bottom, how fast will it be going? Assume it starts from rest.
14. The ballistic pendulum is an interesting device which is used to determine the speed of a bullet. It is a large massive block of wood hanging from a long string. A rifle is fired into the block of wood which then moves. The speed of the bullet can be determined from measuring how high the block of wood rises. Explain how this can be done and why. **Hint:** Let  $v$  be the speed of the bullet which has mass  $m$  and let the block of wood have mass  $M$ . By conservation of momentum  $mv = (m + M)V$  where  $V$  is the speed of the block of wood immediately after the collision. Thus the energy is  $\frac{1}{2}(m + M)V^2$  and this block of wood rises to a height of  $h$ . Now use Problem 12.
15. In the experiment of Problem 14, show the kinetic energy before the collision is greater than the kinetic energy after the collision. Thus linear momentum is conserved but energy is not. Such a collision is called inelastic.
16. There is a popular toy consisting of identical steel balls suspended from strings of equal length as illustrated in the following picture.



The ball at the right is lifted and allowed to swing. When it collides with the other balls, the ball on the left is observed to swing away from the others with the same speed the ball on the right had when it collided. Why does this happen? Why don't two or more of the stationary balls start to move, perhaps at a slower speed? This is an example of an elastic collision because energy is conserved. Of course this could change if you fixed things so the balls would stick to each other.

## 15.7 Systems Of Ordinary Differential Equations

Newton's laws of motion yield a system of ordinary differential equations. Many times the resulting system is so complex that it is impossible to solve in terms of known functions. Therefore, there is a natural question about whether there even exist solutions. This is

a profound question which must ultimately be dealt with. This section contains a fairly general result on existence and uniqueness for systems of ordinary differential equations. First of all, recall from Theorem 10.3.4 on Page 233, if  $f$  is continuous on  $[a, b]$ , then  $f \in R([a, b])$ . Therefore, if  $\mathbf{f}$  is a continuous vector valued function defined on  $[a, b]$  it follows from Definition 15.2.1 on Page 352 that  $\mathbf{f} \in R([a, b])$ . Also from this definition, one can obtain the following simple lemma.

**Lemma 15.7.1** *Let  $\mathbf{f}$  be a continuous function defined on  $[a, b]$  having values in  $\mathbb{R}^n$ . Then  $|\mathbf{f}|$  is also continuous and*

$$\left| \int_a^b \mathbf{f}(t) dt \right| \leq \int_a^b |\mathbf{f}|(t) dt.$$

**Proof:** By the triangle inequality,

$$||\mathbf{f}(t)| - |\mathbf{f}(s)|| \leq |\mathbf{f}(t) - \mathbf{f}(s)| \equiv \left( \sum_{i=1}^n |f_i(t) - f_i(s)|^2 \right)^{1/2}$$

and since  $\mathbf{f}$  is continuous, each  $f_i$  is also and so  $\lim_{s \rightarrow t} |\mathbf{f}(s)| = |\mathbf{f}(t)|$  showing  $\mathbf{f}$  is continuous at  $t$ . It only remains to prove the inequality. To do this, note that if  $\mathbf{u} \in \mathbb{R}^n$ , there exists a vector,  $\mathbf{v}$  such that  $|\mathbf{v}| = 1$  and  $\mathbf{u} \cdot \mathbf{v} = |\mathbf{u}|$ . In fact you may let  $\mathbf{v} = \mathbf{u}/|\mathbf{u}|$  if  $\mathbf{u} \neq \mathbf{0}$ . If  $\mathbf{u} = \mathbf{0}$ , let  $\mathbf{v} = \mathbf{e}_1$ . Therefore, let  $\mathbf{v}$  be such a unit vector with the property that

$$\mathbf{v} \cdot \int_a^b \mathbf{f}(t) dt = \left| \int_a^b \mathbf{f}(t) dt \right|.$$

Then letting  $\mathbf{v} = (v_1, \dots, v_n)$ ,

$$\begin{aligned} \mathbf{v} \cdot \int_a^b \mathbf{f}(t) dt &= v_1 \int_a^b f_1(t) dt + \dots + v_n \int_a^b f_n(t) dt \\ &= \int_a^b \mathbf{v} \cdot \mathbf{f}(t) dt \leq \int_a^b |\mathbf{f}(t)| dt \end{aligned}$$

by the Cauchy Schwarz inequality. This proves the lemma.

### 15.7.1 Picard Iteration

Suppose that  $\mathbf{f} : [a, b] \times \mathbb{R}^n \rightarrow \mathbb{R}^n$  satisfies the following two conditions.

$$|\mathbf{f}(t, \mathbf{x}) - \mathbf{f}(t, \mathbf{x}_1)| \leq K |\mathbf{x} - \mathbf{x}_1|, \quad (15.13)$$

$$\mathbf{f} \text{ is continuous and bounded.} \quad (15.14)$$

The first of these conditions is known as a Lipschitz condition.

**Lemma 15.7.2** *Suppose  $\mathbf{x} : [a, b] \rightarrow \mathbb{R}^n$  is a continuous function and  $c \in [a, b]$ . Then  $\mathbf{x}$  is a solution to the initial value problem,*

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}), \quad \mathbf{x}(c) = \mathbf{x}_0 \quad (15.15)$$

*if and only if  $\mathbf{x}$  is a solution to the integral equation,*

$$\mathbf{x}(t) = \mathbf{x}_0 + \int_c^t \mathbf{f}(s, \mathbf{x}(s)) ds. \quad (15.16)$$

**Proof:** If  $\mathbf{x}$  solves (15.16), then since  $\mathbf{f}$  is continuous, the fundamental theorem of calculus, Theorem 15.2.3 on Page 353 can be used to differentiate both sides and obtain  $\mathbf{x}'(t) = \mathbf{f}(t, \mathbf{x}(t))$ . Also, letting  $t = c$  on both sides, gives  $\mathbf{x}(c) = \mathbf{x}_0$ . Conversely, if  $\mathbf{x}$  is a solution of the initial value problem, integrate both sides from  $c$  to  $t$  obtaining (15.16). This proves the lemma.

It follows from this lemma that the initial value problem, (15.15) is equivalent to the integral equation (15.16) and so it suffices to study this integral equation. The most famous technique for studying this integral equation is the method of Picard iteration. In this method, you begin with an initial function,  $\mathbf{x}_0(t) \equiv \mathbf{x}_0$  and then iterate as follows.

$$\begin{aligned}\mathbf{x}_1(t) &\equiv \mathbf{x}_0 + \int_c^t \mathbf{f}(s, \mathbf{x}_0(s)) ds \\ &= \mathbf{x}_0 + \int_c^t \mathbf{f}(s, \mathbf{x}_0) ds, \\ \mathbf{x}_2(t) &\equiv \mathbf{x}_0 + \int_c^t \mathbf{f}(s, \mathbf{x}_1(s)) ds,\end{aligned}$$

and if  $\mathbf{x}_{k-1}(s)$  has been determined,

$$\mathbf{x}_k(t) \equiv \mathbf{x}_0 + \int_c^t \mathbf{f}(s, \mathbf{x}_{k-1}(s)) ds.$$

These definitions imply some simple estimates. Let

$$M_{\mathbf{f}} \equiv \max \{ |\mathbf{f}(s, \mathbf{x}_0)| : s \in [a, b], \mathbf{x} \in \mathbb{R}^n \}.$$

Then for any  $t \in [a, b]$ ,

$$\begin{aligned}|\mathbf{x}_1(t) - \mathbf{x}_0| &\leq \\ \left| \int_c^t |\mathbf{f}(s, \mathbf{x}_0)| ds \right| &\leq M_{\mathbf{f}} |t - c|. \end{aligned} \tag{15.17}$$

Now using this estimate and the Lipschitz condition for  $\mathbf{f}$ ,

$$\begin{aligned}|\mathbf{x}_2(t) - \mathbf{x}_1(t)| &\leq \left| \int_c^t |\mathbf{f}(s, \mathbf{x}_1(s)) - \mathbf{f}(s, \mathbf{x}_0)| ds \right| \\ &\leq \left| \int_c^t K |\mathbf{x}_1(s) - \mathbf{x}_0| ds \right| \leq \\ K M_{\mathbf{f}} \left| \int_c^t |s - c| ds \right| &\leq K M_{\mathbf{f}} \frac{|t - c|^2}{2}. \end{aligned} \tag{15.18}$$

Continuing in this way leads to the following lemma.

**Lemma 15.7.3** *Let  $k \geq 2$ . Then for any  $t \in [a, b]$ ,*

$$|\mathbf{x}_k(t) - \mathbf{x}_{k-1}(t)| \leq M_{\mathbf{f}} K^{k-1} \frac{|t - c|^k}{k!}. \tag{15.19}$$



**Proof:** The estimate has been verified when  $k = 2$ . Assume it is true for  $k$ . Then the Lipschitz condition on  $\mathbf{f}$  and Lemma 15.7.1 implies

$$\begin{aligned} |\mathbf{x}_{k+1}(t) - \mathbf{x}_k(t)| &= \left| \mathbf{x}_0 + \int_c^t \mathbf{f}(s, \mathbf{x}_k(s)) ds - \left( \mathbf{x}_0 + \int_c^t \mathbf{f}(s, \mathbf{x}_{k+1}(s)) ds \right) \right| \\ &= \left| \int_c^t |\mathbf{f}(s, \mathbf{x}_k(s)) - \mathbf{f}(s, \mathbf{x}_{k+1}(s))| ds \right| \\ &\leq \left| \int_c^t K |\mathbf{x}_k(s) - \mathbf{x}_{k+1}(s)| ds \right| \end{aligned}$$

which, by the induction hypothesis, is dominated by

$$\begin{aligned} &\leq \left| \int_c^t K M_{\mathbf{f}} K^{k-1} \frac{|s-c|^{k-1}}{(k-1)!} ds \right| \\ &\leq M_{\mathbf{f}} K^k \frac{|t-c|^{k+1}}{(k+1)!}. \end{aligned}$$

This proves the lemma.

**Lemma 15.7.4** For each  $t \in [a, b]$ ,  $\{\mathbf{x}_k(t)\}_{k=1}^{\infty}$  is a Cauchy sequence.

**Proof:** Pick such a  $t \in [a, b]$ . Then if  $k > l$ , the triangle inequality and the estimate of the above lemma imply

$$\begin{aligned} |\mathbf{x}_k(t) - \mathbf{x}_l(t)| &\leq \sum_{r=l}^{k-1} |\mathbf{x}_{r+1}(t) - \mathbf{x}_r(t)| \leq \\ M_{\mathbf{f}} \sum_{r=l}^{\infty} K^r \frac{|t-c|^{r+1}}{(r+1)!} &\leq M_{\mathbf{f}} (b-a) \sum_{r=l}^{\infty} \frac{(K(b-a))^r}{(r+1)!} \end{aligned} \quad (15.20)$$

a quantity which converges to zero as  $l \rightarrow \infty$  due to the convergence of the series,  $\sum_{r=0}^{\infty} \frac{(K(b-a))^r}{(r+1)!}$ , a fact which is easily seen by an application of the ratio test. This shows that for every  $\varepsilon > 0$ , there exists  $L$  such that if  $k, l > L$ , then  $|\mathbf{x}_k(t) - \mathbf{x}_l(t)| < \varepsilon$ . In other words,  $\{\mathbf{x}_k(t)\}_{k=1}^{\infty}$  is a Cauchy sequence. This proves the lemma.

Since  $\{\mathbf{x}_k(t)\}_{k=1}^{\infty}$  is a Cauchy sequence, denote by  $\mathbf{x}(t)$  the point to which it converges. Letting  $k \rightarrow \infty$  in (15.20),

$$|\mathbf{x}(t) - \mathbf{x}_l(t)| \leq M_{\mathbf{f}} \sum_{r=l}^{\infty} \frac{(K(b-a))^r}{(r+1)!} < \varepsilon \quad (15.21)$$

whenever  $l$  is sufficiently large. Since the right side of the inequality does not depend on  $t$ , it follows that for all  $\varepsilon > 0$ , there exists  $L$  such that if  $l \geq L$ , then for all  $t \in [a, b]$ ,  $|\mathbf{x}(t) - \mathbf{x}_l(t)| < \varepsilon$ .

**Lemma 15.7.5** The function  $t \rightarrow \mathbf{x}(t)$  is continuous and

$$\lim_{l \rightarrow \infty} (\sup \{|\mathbf{x}(t) - \mathbf{x}_l(t)| : t \in [a, b]\}) = 0. \quad (15.22)$$

**Proof:** Formula (15.22) was established above in (15.21).

Let  $\varepsilon > 0$  be given and let  $t \in [a, b]$ . Then letting  $l$  be large enough that

$$\begin{aligned} \lim_{l \rightarrow \infty} (\sup \{ |\mathbf{x}(t) - \mathbf{x}_l(t)| : t \in [a, b] \}) &< \varepsilon/3, \\ |\mathbf{x}(s) - \mathbf{x}(t)| &\leq \\ |\mathbf{x}(s) - \mathbf{x}_l(s)| + |\mathbf{x}_l(s) - \mathbf{x}_l(t)| + |\mathbf{x}_l(t) - \mathbf{x}(t)| \\ &\leq \frac{2\varepsilon}{3} + |\mathbf{x}_l(s) - \mathbf{x}_l(t)|. \end{aligned}$$

By the continuity of  $\mathbf{x}_l$ , the last term is dominated by  $\frac{\varepsilon}{3}$  whenever  $|s - t|$  is small enough. This verifies continuity of  $t \rightarrow \mathbf{x}(t)$  and proves the lemma.

Letting  $l$  be large enough that

$$\begin{aligned} \lim_{l \rightarrow \infty} (\sup \{ |\mathbf{x}(t) - \mathbf{x}_l(t)| : t \in [a, b] \}) &< \frac{\varepsilon}{K(b-a)}, \\ \left| \int_c^t \mathbf{f}(s, \mathbf{x}_l(s)) ds - \int_c^t \mathbf{f}(s, \mathbf{x}(s)) ds \right| \\ &\leq \left| \int_c^t |\mathbf{f}(s, \mathbf{x}_l(s)) - \mathbf{f}(s, \mathbf{x}(s))| ds \right| \\ &\leq \left| \int_c^t K |\mathbf{x}_l(s) - \mathbf{x}(s)| ds \right| \\ &< K(b-a) \frac{\varepsilon}{K(b-a)} = \varepsilon. \end{aligned}$$

Since  $\varepsilon$  is arbitrary, this verifies that

$$\lim_{l \rightarrow \infty} \int_c^t \mathbf{f}(s, \mathbf{x}_l(s)) ds = \int_c^t \mathbf{f}(s, \mathbf{x}(s)) ds.$$

It follows that

$$\begin{aligned} \mathbf{x}(t) &= \lim_{k \rightarrow \infty} \mathbf{x}_k(t) \\ &= \lim_{k \rightarrow \infty} \left( \mathbf{x}_0 + \int_c^t \mathbf{f}(s, \mathbf{x}_{k-1}(s)) ds \right) \\ &= \mathbf{x}_0 + \int_c^t \mathbf{f}(s, \mathbf{x}(s)) ds \end{aligned}$$

and so by Lemma 15.7.2,  $\mathbf{x}$  is a solution to the initial value problem, (15.15). This proves the existence part of the following theorem.

**Theorem 15.7.6** *Let  $\mathbf{f}$  satisfy (15.13) and (15.14). Then there exists a unique solution to the initial value problem, (15.15).*

**Proof:** It only remains to verify the uniqueness assertion. Therefore, assume  $\mathbf{x}$  and  $\mathbf{x}_1$  are two solutions to the initial value problem. Then by Lemma 15.7.2 it follows that

$$\mathbf{x}(t) - \mathbf{x}_1(t) = \int_c^t (\mathbf{f}(s, \mathbf{x}(s)) - \mathbf{f}(s, \mathbf{x}_1(s))) ds.$$

Suppose first that  $t < c$ . Then this equation implies that for all such  $t$ ,

$$\begin{aligned} |\mathbf{x}(t) - \mathbf{x}_1(t)| &\leq \int_t^c |\mathbf{f}(s, \mathbf{x}(s)) - \mathbf{f}(s, \mathbf{x}_1(s))| \\ &\leq \int_t^c K |\mathbf{x}(s) - \mathbf{x}_1(s)| ds. \end{aligned}$$

Letting  $g(t) = \int_t^c |\mathbf{x}(s) - \mathbf{x}_1(s)| ds$ , it follows that  $-g'(t) = |\mathbf{x}(t) - \mathbf{x}_1(t)|$  and so

$$-g'(t) \leq Kg(t)$$

which implies  $0 \leq g'(t) + Kg(t)$  and consequently,

$$0 \leq (e^{Kt}g(t))'$$

which implies integration from  $t$  to  $c$  gives

$$0 \leq e^{Kc}g(c) - e^{Kt}g(t).$$

But  $g(c) = 0$  and so  $0 \geq g(t) \geq 0$ . Hence  $g(t) = 0$  for all  $t \leq c$  and so  $\mathbf{x}(t) = \mathbf{x}_1(t)$  for these values of  $t$ .

Now suppose that  $t \geq c$ . Then

$$\begin{aligned} |\mathbf{x}(t) - \mathbf{x}_1(t)| &\leq \int_c^t |\mathbf{f}(s, \mathbf{x}(s)) - \mathbf{f}(s, \mathbf{x}_1(s))| \\ &\leq K \int_c^t |\mathbf{x}(s) - \mathbf{x}_1(s)| ds. \end{aligned}$$

Letting  $h(t) = \int_c^t |\mathbf{x}(s) - \mathbf{x}_1(s)| ds$ , it follows that

$$h'(t) \leq Kh(t)$$

and so

$$(e^{-Kt}h(t))' \leq 0$$

which means  $e^{-Kt}h(t) - e^{-Kc}h(c) \leq 0$ . But  $h(c) = 0$  and so this requires

$$h(t) \equiv \int_c^t |\mathbf{x}(s) - \mathbf{x}_1(s)| ds \leq 0$$

and so  $\mathbf{x}(t) = \mathbf{x}_1(t)$  whenever  $t \geq c$ . Therefore,  $\mathbf{x}_1 = \mathbf{x}$ . This proves uniqueness and establishes the theorem.

## 15.7.2 Numerical Methods For Differential Equations

The above considers a fundamental philosophical question about existence of solutions to the initial value problem. In principle, it yields a method for finding the solution as well but it is never used this way. There are much simpler and better ways to obtain approximate solutions to differential equations than using Picard iteration. The most primitive way of doing this is the Euler method. Suppose you want to find the solution to the initial value problem,

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}), \quad \mathbf{x}(0) = \mathbf{x}_0,$$

for  $t \in [0, T]$ , where  $\mathbf{f}$  satisfies the conditions given above for existence. You split the time interval into  $n$  equal pieces,  $0, \frac{T}{n}, \frac{2T}{n}, \dots, \frac{(n-1)T}{n}, \frac{nT}{n}$ . Letting  $h = \frac{T}{n}$  and  $t_i = ih$ , you then replace the derivative with a difference quotient as follows.

$$\frac{\mathbf{x}(t_i + h) - \mathbf{x}(t_i)}{h} = \mathbf{f}(t_i, \mathbf{x}(t_i)), \quad \mathbf{x}(0) = \mathbf{x}_0. \quad (15.23)$$

Noting that  $t_i + h = t_{i+1}$ , this can be written more simply as

$$\mathbf{x}(t_{i+1}) = \mathbf{x}(t_i) + h\mathbf{f}(t_i, \mathbf{x}(t_i)), \quad \mathbf{x}(0) = \mathbf{x}_0. \quad (15.24)$$

Now this can be used to find an approximate solution as follows.  $\mathbf{x}_0$  is known because it is given. From (15.24)  $\mathbf{x}(t_1) = \mathbf{x}_0 + h\mathbf{f}(0, \mathbf{x}_0)$ . Having found  $\mathbf{x}(t_1)$ , this can be substituted back in to (15.24) to find  $\mathbf{x}(t_2) = \mathbf{x}(t_1) + h\mathbf{f}(t_1, \mathbf{x}(t_1))$ . Using  $\mathbf{x}(t_2)$  one used (15.24) to find  $\mathbf{x}(t_3)$  and you continue doing this till you get  $\mathbf{x}(t_n)$ . For  $t \in (t_k, t_{k+1})$  you define

$$\mathbf{x}(t) \equiv \frac{(t_{k+1} - t)}{h}\mathbf{x}(t_k) + \frac{(t - t_k)}{h}\mathbf{x}(t_{k+1}).$$

This whole idea is based on approximating functions with their linear approximations locally. Therefore, the resulting function will be piecewise linear and hopefully will not be too far off. It is customary to write  $\mathbf{x}_i$  instead of  $\mathbf{x}(t_i)$ .

How well does it work? In general, when you are working with numerical methods you test them on known examples and see how well they perform. In the following example, Euler's method will be applied to an initial value problem whose solution is known.

**Example 15.7.7** Solve the initial value problem  $y' = y$ ,  $y(0) = 1$  for  $t \in [0, 5]$  using Euler's method with a step size of .01 and a step size of .001.

The correct answer is  $y = e^x$ . Thus  $y(5) = e^5 = 148.413\,159\,102\,576\,6$ . Now from Euler's method,

$$y_1 = 1 + .01(0) = 1$$

Now

$$y_2 = 1 + (.01)1 = 1.01.$$

Next,

$$y_3 = 1.01 + (.01)1.01 = (1.01)^2.$$

Continuing this way, and using the pattern which is emerging,

$$y_{500} = (1.01)^{500} = 144.77$$

and this is the value which Euler's method gives for  $y(5)$ . If you used the step size .001, you would have  $(1.001)^{5000} = 148.042$  which is pretty close.

Of course this was an easy example because it led to a pattern which was easy to apply. In general, there will be no pattern and you just have to iterate things till you get to the end. This is done by a computer.

You need a much better method than the Euler method if you want to solve ordinary differential equations. There are many methods and the detailed study of them is outside the scope of this book. One of the most famous is the Runge Kutta fourth order method. To solve  $\mathbf{x}' = \mathbf{f}(t, \mathbf{x})$ ,  $\mathbf{x}(0) = \mathbf{x}_0$ , you do the following to go from  $t_n$  to  $t_{n+1}$ .

$$\begin{aligned} \mathbf{k}_1 &\equiv h\mathbf{f}(t_n, \mathbf{x}_n) \\ \mathbf{k}_2 &\equiv h\mathbf{f}\left(t_n + \frac{1}{2}h, \mathbf{x}_n + \frac{1}{2}\mathbf{k}_1\right) \\ \mathbf{k}_3 &\equiv h\mathbf{f}\left(t_n + \frac{1}{2}h, \mathbf{x}_n + \frac{1}{2}\mathbf{k}_2\right) \\ \mathbf{k}_4 &\equiv h\mathbf{f}(t_n + h, \mathbf{x}_n + \mathbf{k}_3). \end{aligned}$$

Then

$$\mathbf{x}_{n+1} \equiv \mathbf{x}_n + \frac{1}{6} (\mathbf{k}_1 + 2\mathbf{k}_2 + 2\mathbf{k}_3 + \mathbf{k}_4).$$

This method is called fourth order because it can be shown that the difference between the true solution and this approximate solution is dominated by  $Ch^4$  where  $C$  is some constant independent of the choice of  $h$ .

Even better schemes have been devised. The computer algebra system, Maple, uses one such system which is a fifth order method. The important thing for you to realize at this point is that the initial value problem for an ordinary differential equation is a significant problem and it has been solved quite well.

## 15.8 Exercises

1. Consider the model problem  $y' = y$ ,  $y(0) = 1$  on the interval  $[0, 1]$ . Solve it numerically for  $h = .2$  using the Euler method and compare with the true solution. Now use the Runge Kutta method for  $h = .2$  and compare with the true solution.
2. Show that if  $y' = f(x)$ ,  $y(0) = 0$  so the differential equation does not depend on  $y$ , then the Runge Kutta method above reduces to Simpson's rule for numerical integration.
3. Show that in going from 0 to  $t_1$ , the difference between the Euler solution and the true solution is less than a constant times  $h^2$ . Assume anything you like about the smoothness of the function,  $f$ . Now show how these errors add up and conclude the Euler method is of order 1. Thus the difference between the true solution and the Euler solution is less than  $Ch$ .
4. Work the above model problem using the following method known as a predictor corrector method.

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \frac{h}{2} [\mathbf{f}(t_n, \mathbf{x}_n) + \mathbf{f}(t_{n+1}, \mathbf{x}_{n+1}^*)]$$

where here  $\mathbf{x}_{n+1}^* = \mathbf{x}_n + h\mathbf{f}(t_n, \mathbf{x}_n)$ . This  $\mathbf{x}_{n+1}^*$  is the prediction. Then  $\mathbf{x}_{n+1}$  is the correction. Use  $h = .2$  and compare with the true solution.



# Line Integrals

The concept of the integral can be extended to functions which are not defined on an interval of the real line but on some curve in  $\mathbb{R}^n$ . This is done by defining things in such a way that the more general concept reduces to the earlier notion. First it is necessary to consider what is meant by arc length.

## 16.1 Arc Length And Orientations

The application of the integral considered here is the concept of the length of a curve.  $C$  is a smooth curve in  $\mathbb{R}^n$  if there exists an interval,  $[a, b] \subseteq \mathbb{R}$  and functions  $x_i : [a, b] \rightarrow \mathbb{R}$  such that the following conditions hold

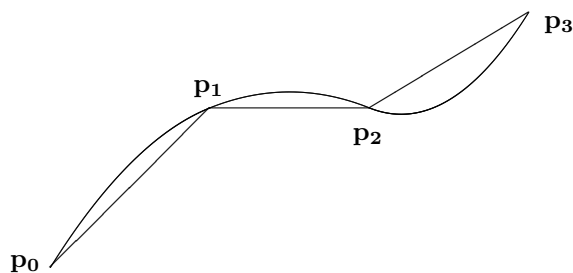
1.  $x_i$  is continuous on  $[a, b]$ .
2.  $x'_i$  exists and is continuous and bounded on  $[a, b]$ , with  $x'_i(a)$  defined as the derivative from the right,

$$\lim_{h \rightarrow 0+} \frac{x_i(a+h) - x_i(a)}{h},$$

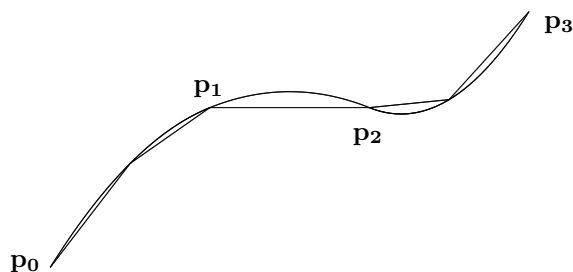
and  $x'_i(b)$  defined similarly as the derivative from the left.

3. For  $\mathbf{p}(t) \equiv (x_1(t), \dots, x_n(t))$ ,  $t \rightarrow \mathbf{p}(t)$  is one to one on  $(a, b)$ .
4.  $|\mathbf{p}'(t)| \equiv \left( \sum_{i=1}^n |x'_i(t)|^2 \right)^{1/2} \neq 0$  for all  $t \in [a, b]$ .
5.  $C = \cup \{(x_1(t), \dots, x_n(t)) : t \in [a, b]\}$ .

The functions,  $x_i(t)$ , defined above are giving the coordinates of a point in  $\mathbb{R}^n$  and the list of these functions is called a parameterization for the smooth curve. Note the natural direction of the interval also gives a direction for moving along the curve. Such a direction is called an orientation. The integral is used to define what is meant by the length of such a smooth curve. The earlier treatment in terms of initial value problems will be made more rigorous here. Consider such a smooth curve having parameterization  $(x_1, \dots, x_n)$ . Forming a partition of  $[a, b]$ ,  $a = t_0 < \dots < t_n = b$  and letting  $\mathbf{p}_i = (x_1(t_i), \dots, x_n(t_i))$ , you could consider the polygon formed by lines from  $\mathbf{p}_0$  to  $\mathbf{p}_1$  and from  $\mathbf{p}_1$  to  $\mathbf{p}_2$  and from  $\mathbf{p}_3$  to  $\mathbf{p}_4$  etc. to be an approximation to the curve,  $C$ . The following picture illustrates what is meant by this.



Now consider what happens when the partition is refined by including more points. You can see from the following picture that the polygonal approximation would appear to be even better and that as more points are added in the partition, the sum of the lengths of the line segments seems to get close to something which deserves to be defined as the length of the curve,  $C$ .



Thus the length of the curve is approximated by

$$\sum_{k=1}^n |\mathbf{p}(t_k) - \mathbf{p}(t_{k-1})|.$$

Since the functions in the parameterization are differentiable, it is reasonable to expect this to be close to

$$\sum_{k=1}^n |\mathbf{p}'(t_{k-1})| (t_k - t_{k-1})$$

which is seen to be a Riemann sum for the integral

$$\int_a^b |\mathbf{p}'(t)| dt$$

and it is this integral which is defined as the length of the curve.

Would the same length be obtained if another parameterization were used? This is a very important question because the length of the curve should depend only on the curve itself and not on the method used to trace out the curve. The answer to this question is that the length of the curve does not depend on parameterization. The proof is somewhat technical so is given in the last section of this chapter.

Does the definition of length given above correspond to the usual definition of length in the case when the curve is a line segment? It is easy to see that it does so by considering two points in  $\mathbb{R}^n$ ,  $\mathbf{p}$  and  $\mathbf{q}$ . A parameterization for the line segment joining these two points is

$$f_i(t) \equiv tp_i + (1-t)q_i, \quad t \in [0, 1].$$



Using the definition of length of a smooth curve just given, the length according to this definition is

$$\int_0^1 \left( \sum_{i=1}^n (p_i - q_i)^2 \right)^{1/2} dt = |\mathbf{p} - \mathbf{q}|.$$

Thus this new definition which is valid for smooth curves which may not be straight line segments gives the usual length for straight line segments.

The proof that curve length is well defined for a smooth curve contains a result which deserves to be stated as a corollary. It is proved in Lemma 16.6.2 on Page 385 but the proof is mathematically fairly advanced so it is presented later.

**Corollary 16.1.1** *Let  $C$  be a smooth curve and let  $\mathbf{f} : [a, b] \rightarrow C$  and  $\mathbf{g} : [c, d] \rightarrow C$  be two parameterizations satisfying 1 - 5. Then  $\mathbf{g}^{-1} \circ \mathbf{f}$  is either strictly increasing or strictly decreasing.*

**Definition 16.1.2** *If  $\mathbf{g}^{-1} \circ \mathbf{f}$  is increasing, then  $\mathbf{f}$  and  $\mathbf{g}$  are said to be equivalent parameterizations and this is written as  $\mathbf{f} \sim \mathbf{g}$ . It is also said that the two parameterizations give the same orientation for the curve when  $\mathbf{f} \sim \mathbf{g}$ .*

When the parameterizations are equivalent, they preserve the direction, of motion along the curve and this also shows there are exactly two orientations of the curve since either  $\mathbf{g}^{-1} \circ \mathbf{f}$  is increasing or it is decreasing. This is not hard to believe. In simple language, the message is that there are exactly two directions of motion along a curve.

**Lemma 16.1.3** *The following hold for  $\sim$ .*

$$\mathbf{f} \sim \mathbf{f}, \quad (16.1)$$

$$\text{If } \mathbf{f} \sim \mathbf{g} \text{ then } \mathbf{g} \sim \mathbf{f}, \quad (16.2)$$

$$\text{If } \mathbf{f} \sim \mathbf{g} \text{ and } \mathbf{g} \sim \mathbf{h}, \text{ then } \mathbf{f} \sim \mathbf{h}. \quad (16.3)$$

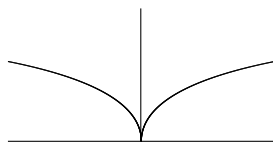
**Proof:** Formula (16.1) is obvious because  $\mathbf{f}^{-1} \circ \mathbf{f}(t) = t$  so it is clearly an increasing function. If  $\mathbf{f} \sim \mathbf{g}$  then  $\mathbf{f}^{-1} \circ \mathbf{g}$  is increasing. Now  $\mathbf{g}^{-1} \circ \mathbf{f}$  must also be increasing because it is the inverse of  $\mathbf{f}^{-1} \circ \mathbf{g}$ . This verifies (16.2). To see (16.3),  $\mathbf{f}^{-1} \circ \mathbf{h} = (\mathbf{f}^{-1} \circ \mathbf{g}) \circ (\mathbf{g}^{-1} \circ \mathbf{h})$  and so since both of these functions are increasing, it follows  $\mathbf{f}^{-1} \circ \mathbf{h}$  is also increasing. This proves the lemma.

The symbol,  $\sim$  is called an equivalence relation. If  $C$  is such a smooth curve just described, and if  $\mathbf{f} : [a, b] \rightarrow C$  is a parameterization of  $C$ , consider  $\mathbf{g}(t) \equiv \mathbf{f}((a+b)-t)$ , also a parameterization of  $C$ . Now by Corollary 16.1.1, if  $\mathbf{h}$  is a parameterization, then if  $\mathbf{f}^{-1} \circ \mathbf{h}$  is not increasing, it must be the case that  $\mathbf{g}^{-1} \circ \mathbf{h}$  is increasing. Consequently, either  $\mathbf{h} \sim \mathbf{g}$  or  $\mathbf{h} \sim \mathbf{f}$ . These parameterizations,  $\mathbf{h}$ , which satisfy  $\mathbf{h} \sim \mathbf{f}$  are called the equivalence class determined by  $\mathbf{f}$  and those  $\mathbf{h} \sim \mathbf{g}$  are called the equivalence class determined by  $\mathbf{g}$ . These two classes are called orientations of  $C$ . They give the direction of motion on  $C$ . You see that going from  $\mathbf{f}$  to  $\mathbf{g}$  corresponds to tracing out the curve in the opposite direction.

Sometimes people wonder why it is required, in the definition of a smooth curve that  $\mathbf{p}'(t) \neq \mathbf{0}$ . Imagine  $t$  is time and  $\mathbf{p}(t)$  gives the location of a point in space. If  $\mathbf{p}'(t)$  is allowed to equal zero, the point can stop and change directions abruptly, producing a pointy place in  $C$ . Here is an example.

**Example 16.1.4** *Graph the curve  $(t^3, t^2)$  for  $t \in [-1, 1]$ .*

In this case,  $t = x^{1/3}$  and so  $y = x^{2/3}$ . Thus the graph of this curve looks like the picture below. Note the pointy place. Such a curve should not be considered smooth.

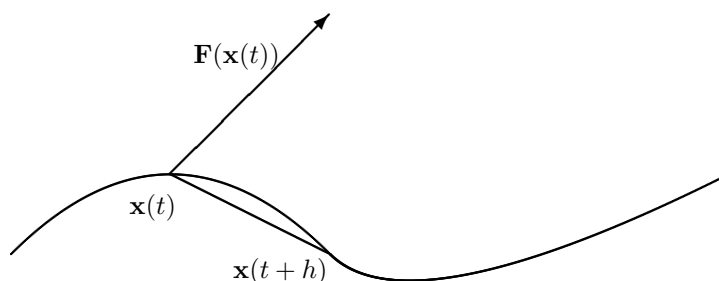


## 16.2 Line Integrals And Work

Let  $C$  be a smooth curve contained in  $\mathbb{R}^p$ . A curve,  $C$  is an “oriented curve” if the only parameterizations considered are those which lie in exactly one of the two equivalence classes, each of which is called an “orientation”. In simple language, orientation specifies a direction over which motion along the curve is to take place. Thus, it specifies the order in which the points of  $C$  are encountered. The pair of concepts consisting of the set of points making up the curve along with a direction of motion along the curve is called an oriented curve.

**Definition 16.2.1** Suppose  $\mathbf{F}(\mathbf{x}) \in \mathbb{R}^p$  is given for each  $\mathbf{x} \in C$  where  $C$  is a smooth oriented curve and suppose  $\mathbf{x} \rightarrow \mathbf{F}(\mathbf{x})$  is continuous. The mapping  $\mathbf{x} \rightarrow \mathbf{F}(\mathbf{x})$  is called a vector field. In the case that  $\mathbf{F}(\mathbf{x})$  is a force, it is called a force field.

Next the concept of work done by a force field,  $\mathbf{F}$  on an object as it moves along the curve,  $C$ , in the direction determined by the given orientation of the curve will be defined. This is new. Earlier the work done by a force which acts on an object moving in a straight line was discussed but here the object moves over a curve. In order to define what is meant by the work, consider the following picture.



In this picture, the work done by a constant force,  $\mathbf{F}$  on an object which moves from the point  $\mathbf{x}(t)$  to the point  $\mathbf{x}(t+h)$  along the straight line shown would equal  $\mathbf{F} \cdot (\mathbf{x}(t+h) - \mathbf{x}(t))$ . It is reasonable to assume this would be a good approximation to the work done in moving along the curve joining  $\mathbf{x}(t)$  and  $\mathbf{x}(t+h)$  provided  $h$  is small enough. Also, provided  $h$  is small,

$$\mathbf{x}(t+h) - \mathbf{x}(t) \approx \mathbf{x}'(t)h$$

where the wriggly equal sign indicates the two quantities are close. In the notation of Leibniz, one writes  $dt$  for  $h$  and

$$dW = \mathbf{F}(\mathbf{x}(t)) \cdot \mathbf{x}'(t) dt$$

or in other words,

$$\frac{dW}{dt} = \mathbf{F}(\mathbf{x}(t)) \cdot \mathbf{x}'(t).$$

Defining the total work done by the force at  $t = 0$ , corresponding to the first endpoint of the curve, to equal zero, the work would satisfy the following initial value problem.

$$\frac{dW}{dt} = \mathbf{F}(\mathbf{x}(t)) \cdot \mathbf{x}'(t), \quad W(a) = 0.$$

This motivates the following definition of work.

**Definition 16.2.2** Let  $\mathbf{F}(\mathbf{x})$  be given above. Then the work done by this force field on an object moving over the curve  $C$  in the direction determined by the specified orientation is defined as

$$\int_C \mathbf{F} \cdot d\mathbf{R} \equiv \int_a^b \mathbf{F}(\mathbf{x}(t)) \cdot \mathbf{x}'(t) dt$$

where the function,  $\mathbf{x}$  is one of the allowed parameterizations of  $C$  in the given orientation of  $C$ . In other words, there is an interval,  $[a, b]$  and as  $t$  goes from  $a$  to  $b$ ,  $\mathbf{x}(t)$  moves in the direction determined from the given orientation of the curve.

**Theorem 16.2.3** The symbol,  $\int_C \mathbf{F} \cdot d\mathbf{R}$ , is well defined in the sense that every parameterization in the given orientation of  $C$  gives the same value for  $\int_C \mathbf{F} \cdot d\mathbf{R}$ .

**Proof:** Suppose  $\mathbf{g} : [c, d] \rightarrow C$  is another allowed parameterization. Thus  $\mathbf{g}^{-1} \circ \mathbf{f}$  is an increasing function,  $\phi$ . Then since  $\phi$  is increasing,

$$\begin{aligned} \int_c^d \mathbf{F}(\mathbf{g}(s)) \cdot \mathbf{g}'(s) ds &= \int_a^b \mathbf{F}(\mathbf{g}(\phi(t))) \cdot \mathbf{g}'(\phi(t)) \phi'(t) dt \\ &= \int_a^b \mathbf{F}(\mathbf{f}(t)) \cdot \frac{d}{dt} (\mathbf{g}(\mathbf{g}^{-1} \circ \mathbf{f}(t))) dt = \int_a^b \mathbf{F}(\mathbf{f}(t)) \cdot \mathbf{f}'(t) dt. \end{aligned}$$

This proves the theorem.

Regardless the physical interpretation of  $\mathbf{F}$ , this is called the line integral. When  $\mathbf{F}$  is interpreted as a force, the line integral measures the extent to which the motion over the curve in the indicated direction is aided by the force. If the net effect of the force on the object is to impede rather than to aid the motion, this will show up as the work being negative.

Does the concept of work as defined here coincide with the earlier concept of work when the object moves over a straight line when acted on by a constant force?

Let  $\mathbf{p}$  and  $\mathbf{q}$  be two points in  $\mathbb{R}^n$  and suppose  $\mathbf{F}$  is a constant force acting on an object which moves from  $\mathbf{p}$  to  $\mathbf{q}$  along the straight line joining these points. Then the work done is  $\mathbf{F} \cdot (\mathbf{q} - \mathbf{p})$ . Is the same thing obtained from the above definition? Let  $\mathbf{x}(t) \equiv \mathbf{p} + t(\mathbf{q} - \mathbf{p})$ ,  $t \in [0, 1]$  be a parameterization for this oriented curve, the straight line in the direction from  $\mathbf{p}$  to  $\mathbf{q}$ . Then  $\mathbf{x}'(t) = \mathbf{q} - \mathbf{p}$  and  $\mathbf{F}(\mathbf{x}(t)) = \mathbf{F}$ . Therefore, the above definition yields

$$\int_0^1 \mathbf{F} \cdot (\mathbf{q} - \mathbf{p}) dt = \mathbf{F} \cdot (\mathbf{q} - \mathbf{p}).$$

Therefore, the new definition adds to but does not contradict the old one.

**Example 16.2.4** Suppose for  $t \in [0, \pi]$  the position of an object is given by  $\mathbf{r}(t) = t\mathbf{i} + \cos(2t)\mathbf{j} + \sin(2t)\mathbf{k}$ . Also suppose there is a force field defined on  $\mathbb{R}^3$ ,  $\mathbf{F}(x, y, z) \equiv 2xy\mathbf{i} + x^2\mathbf{j} + \mathbf{k}$ . Find

$$\int_C \mathbf{F} \cdot d\mathbf{R}$$

where  $C$  is the curve traced out by this object which has the orientation determined by the direction of increasing  $t$ .

To find this line integral use the above definition and write

$$\int_C \mathbf{F} \cdot d\mathbf{R} = \int_0^\pi (2t(\cos(2t)), t^2, 1) \cdot (1, -2\sin(2t), 2\cos(2t)) dt$$

In evaluating this replace the  $x$  in the formula for  $\mathbf{F}$  with  $t$ , the  $y$  in the formula for  $\mathbf{F}$  with  $\cos(2t)$  and the  $z$  in the formula for  $\mathbf{F}$  with  $\sin(2t)$  because these are the values of these variables which correspond to the value of  $t$ . Taking the dot product, this equals the following integral.

$$\int_0^\pi (2t \cos 2t - 2(\sin 2t)t^2 + 2 \cos 2t) dt = \pi^2$$

## 16.3 Exercises

- Suppose for  $t \in [0, 2\pi]$  the position of an object is given by  $\mathbf{r}(t) = t\mathbf{i} + \cos(2t)\mathbf{j} + \sin(2t)\mathbf{k}$ . Also suppose there is a force field defined on  $\mathbb{R}^3$ ,  $\mathbf{F}(x, y, z) \equiv 2xy\mathbf{i} + (x^2 + 2zy)\mathbf{j} + y^2\mathbf{k}$ . Find

$$\int_C \mathbf{F} \cdot d\mathbf{R}$$

where  $C$  is the curve traced out by this object which has the orientation determined by the direction of increasing  $t$ .

- Here is a vector field,  $(y, x + z^2, 2yz)$  and here is the parameterization of a curve,  $C$ .  $\mathbf{R}(t) = (\cos 2t, 2\sin 2t, t)$  where  $t$  goes from 0 to  $\pi/4$ . Find  $\int_C \mathbf{F} \cdot d\mathbf{R}$ .
- If  $f$  and  $g$  are both increasing functions, show  $f \circ g$  is an increasing function also. Assume anything you like about the domains of the functions.
- Suppose for  $t \in [0, 3]$  the position of an object is given by  $\mathbf{r}(t) = t\mathbf{i} + t\mathbf{j} + t\mathbf{k}$ . Also suppose there is a force field defined on  $\mathbb{R}^3$ ,  $\mathbf{F}(x, y, z) \equiv yz\mathbf{i} + xz\mathbf{j} + xy\mathbf{k}$ . Find

$$\int_C \mathbf{F} \cdot d\mathbf{R}$$

where  $C$  is the curve traced out by this object which has the orientation determined by the direction of increasing  $t$ . Repeat the problem for  $\mathbf{r}(t) = t\mathbf{i} + t^2\mathbf{j} + t\mathbf{k}$ .

- Suppose for  $t \in [0, 1]$  the position of an object is given by  $\mathbf{r}(t) = t\mathbf{i} + t\mathbf{j} + t\mathbf{k}$ . Also suppose there is a force field defined on  $\mathbb{R}^3$ ,  $\mathbf{F}(x, y, z) \equiv z\mathbf{i} + xz\mathbf{j} + xy\mathbf{k}$ . Find

$$\int_C \mathbf{F} \cdot d\mathbf{R}$$

where  $C$  is the curve traced out by this object which has the orientation determined by the direction of increasing  $t$ . Repeat the problem for  $\mathbf{r}(t) = t\mathbf{i} + t^2\mathbf{j} + t\mathbf{k}$ .

6. Let  $\mathbf{F}(x, y, z)$  be a given force field and suppose it acts on an object having mass,  $m$  on a curve with parameterization,  $(x(t), y(t), z(t))$  for  $t \in [a, b]$ . Show directly that the work done equals the difference in the kinetic energy. **Hint:**

$$\int_a^b \mathbf{F}(x(t), y(t), z(t)) \cdot (x'(t), y'(t), z'(t)) dt =$$

$$\int_a^b m(x''(t), y''(t), z''(t)) \cdot (x'(t), y'(t), z'(t)) dt,$$

etc.

## 16.4 Motion On A Space Curve

A fly buzzing around the room, a person riding a roller coaster, and a satellite orbiting the earth all have something in common. They are moving over some sort of curve in three dimensions.

Denote by  $\mathbf{R}(s)$  the function which takes  $s$  to a point on this curve where  $s$  is arc length. Thus  $\mathbf{R}(s)$  equals the point on the curve which occurs when you have traveled a distance of  $s$  along the curve from one end. This is known as the parameterization of the curve in terms of arc length. Note also that it incorporates an orientation on the curve because there are exactly two ends you could begin measuring length from. In this section, assume anything about smoothness and continuity to make the following manipulations valid. In particular, assume that  $\mathbf{R}'$  exists and is continuous.

**Lemma 16.4.1** Define  $\mathbf{T}(s) \equiv \mathbf{R}'(s)$ . Then  $|\mathbf{T}(s)| = 1$  and if  $\mathbf{T}'(s) \neq 0$ , then there exists a unit vector,  $\mathbf{N}(s)$  perpendicular to  $\mathbf{T}(s)$  and a scalar valued function,  $\kappa(s)$  with  $\mathbf{T}'(s) = \kappa(s)\mathbf{N}(s)$ .

**Proof:** First,  $s = \int_0^s |\mathbf{R}'(r)| dr$  because of the definition of arc length. Therefore, from the fundamental theorem of calculus,  $1 = |\mathbf{R}'(s)| = |\mathbf{T}(s)|$ . Therefore,  $\mathbf{T} \cdot \mathbf{T} = 1$  and so upon differentiating this on both sides, yields  $\mathbf{T}' \cdot \mathbf{T} + \mathbf{T} \cdot \mathbf{T}' = 0$  which shows  $\mathbf{T} \cdot \mathbf{T}' = 0$ . Therefore, the vector,  $\mathbf{T}'$  is perpendicular to the vector,  $\mathbf{T}$ . In case  $\mathbf{T}'(s) \neq 0$ , let  $\mathbf{N}(s) = \frac{\mathbf{T}'(s)}{|\mathbf{T}'(s)|}$  and so  $\mathbf{T}'(s) = |\mathbf{T}'(s)|\mathbf{N}(s)$ , showing the scalar valued function is  $\kappa(s) = |\mathbf{T}'(s)|$ . This proves the lemma.

The radius of curvature is defined as  $\rho = \frac{1}{\kappa}$ . Thus at points where there is a lot of curvature, the radius of curvature is small and at points where the curvature is small, the radius of curvature is large. The plane determined by the two vectors,  $\mathbf{T}$  and  $\mathbf{N}$  is called the osculating plane. It identifies a particular plane which is in a sense tangent to this space curve. In the case where  $|\mathbf{T}'(s)| = 0$  near the point of interest,  $\mathbf{T}(s)$  equals a constant and so the space curve is a straight line which it would be supposed has no curvature. Also, the principal normal is undefined in this case. This makes sense because if there is no curving going on, there is no special direction normal to the curve at such points which could be distinguished from any other direction normal to the curve. In the case where  $|\mathbf{T}'(s)| = 0$ ,  $\kappa(s) = 0$  and the radius of curvature would be considered infinite.

**Definition 16.4.2** The vector,  $\mathbf{T}(s)$  is called the unit tangent vector and the vector,  $\mathbf{N}(s)$  is called the principal normal. The function,  $\kappa(s)$  in the above lemma is called the curvature. When  $\mathbf{T}'(s) \neq 0$  so the principal normal is defined, the vector,  $\mathbf{B}(s) \equiv \mathbf{T}(s) \times \mathbf{N}(s)$  is called the binormal.

The binormal is normal to the osculating plane and  $\mathbf{B}'$  tells how fast this vector changes. Thus it measures the rate at which the curve twists.

**Lemma 16.4.3** *Let  $\mathbf{R}(s)$  be a parameterization of a space curve with respect to arc length and let the vectors,  $\mathbf{T}$ ,  $\mathbf{N}$ , and  $\mathbf{B}$  be as defined above. Then  $\mathbf{B}' = \mathbf{T} \times \mathbf{N}'$  and there exists a scalar function,  $\tau(s)$  such that  $\mathbf{B}' = \tau\mathbf{N}$ .*

**Proof:** From the definition of  $\mathbf{B} = \mathbf{T} \times \mathbf{N}$ , and you can differentiate both sides and get  $\mathbf{B}' = \mathbf{T}' \times \mathbf{N} + \mathbf{T} \times \mathbf{N}'$ . Now recall that  $\mathbf{T}'$  is a multiple called curvature multiplied by  $\mathbf{N}$  so the vectors,  $\mathbf{T}'$  and  $\mathbf{N}$  have the same direction and  $\mathbf{B}' = \mathbf{T} \times \mathbf{N}'$ . Therefore,  $\mathbf{B}'$  is either zero or is perpendicular to  $\mathbf{T}$ . But also, from the definition of  $\mathbf{B}$ ,  $\mathbf{B}$  is a unit vector and so  $\mathbf{B}(s) \cdot \mathbf{B}(s) = 1$ . Differentiating this,  $\mathbf{B}'(s) \cdot \mathbf{B}(s) + \mathbf{B}(s) \cdot \mathbf{B}'(s) = 0$  showing that  $\mathbf{B}'$  is perpendicular to  $\mathbf{B}$  also. Therefore,  $\mathbf{B}'$  is a vector which is perpendicular to both vectors,  $\mathbf{T}$  and  $\mathbf{B}$  and since this is in three dimensions,  $\mathbf{B}'$  must be some scalar multiple of  $\mathbf{N}$  and it is this multiple called  $\tau$ . Thus  $\mathbf{B}' = \tau\mathbf{N}$  as claimed.

Lets go over this last claim a little more. Later it will be done algebraically but for now, the following situation is obtained. There are two vectors,  $\mathbf{T}$  and  $\mathbf{B}$  which are perpendicular to each other and both  $\mathbf{B}'$  and  $\mathbf{N}$  are perpendicular to these two vectors, hence perpendicular to the plane determined by them. Therefore,  $\mathbf{B}'$  must be a multiple of  $\mathbf{N}$ . Take a piece of paper, draw two unit vectors on it which are perpendicular. Then you can see that any two vectors which are perpendicular to this plane must be multiples of each other.

The scalar function,  $\tau$  is called the torsion. In case  $\mathbf{T}' = 0$ , none of this is defined because in this case there is not a well defined osculating plane. The conclusion of the following theorem is called the Serret Frenet formulas.

**Theorem 16.4.4 (Serret Frenet)** *Let  $\mathbf{R}(s)$  be the parameterization with respect to arc length of a space curve and  $\mathbf{T}(s) = \mathbf{R}'(s)$  is the unit tangent vector. Suppose  $|\mathbf{T}'(s)| \neq 0$  so the principal normal,  $\mathbf{N}(s) = \frac{\mathbf{T}'(s)}{|\mathbf{T}'(s)|}$  is defined. The binormal is the vector  $\mathbf{B} \equiv \mathbf{T} \times \mathbf{N}$  so  $\mathbf{T}, \mathbf{N}, \mathbf{B}$  forms a right handed system of unit vectors each of which is perpendicular to every other. Then the following system of differential equations holds in  $\mathbb{R}^3$ .*

$$\mathbf{B}' = \tau\mathbf{N}, \quad \mathbf{T}' = \kappa\mathbf{N}, \quad \mathbf{N}' = -\kappa\mathbf{T} - \tau\mathbf{B}$$

where  $\kappa$  is the curvature and is nonnegative and  $\tau$  is the torsion.

**Proof:**  $\kappa \geq 0$  because  $\kappa = |\mathbf{T}'(s)|$ . The first two equations are already established. To get the third, note that  $\mathbf{B} \times \mathbf{T} = \mathbf{N}$  which follows because  $\mathbf{T}, \mathbf{N}, \mathbf{B}$  is given to form a right handed system of unit vectors each perpendicular to the others. (Use your right hand.) Now take the derivative of this expression. thus

$$\begin{aligned} \mathbf{N}' &= \mathbf{B}' \times \mathbf{T} + \mathbf{B} \times \mathbf{T}' \\ &= \tau\mathbf{N} \times \mathbf{T} + \kappa\mathbf{B} \times \mathbf{N}. \end{aligned}$$

Now recall again that  $\mathbf{T}, \mathbf{N}, \mathbf{B}$  is a right hand system. Thus  $\mathbf{N} \times \mathbf{T} = -\mathbf{B}$  and  $\mathbf{B} \times \mathbf{N} = -\mathbf{T}$ . This establishes the Frenet Serret formulas.

This is an important example of a system of differential equations in  $\mathbb{R}^3$ . It is a remarkable result because it says that from knowledge of the two scalar functions,  $\tau$  and  $\kappa$ , and initial values for  $\mathbf{B}$ ,  $\mathbf{T}$ , and  $\mathbf{N}$  when  $s = 0$  you can obtain the binormal, unit tangent, and principal normal vectors. It is just the solution of an initial value problem of the sort discussed earlier. Having done this, you can reconstruct the entire space curve starting at some point,  $\mathbf{R}_0$  because  $\mathbf{R}'(s) = \mathbf{T}(s)$  and so  $\mathbf{R}(s) = \mathbf{R}_0 + \int_0^s \mathbf{T}'(r) dr$ .

The vectors,  $\mathbf{B}$ ,  $\mathbf{T}$ , and  $\mathbf{N}$  are vectors which are functions of position on the space curve. Often, especially in applications, you deal with a space curve which is parameterized by a function of  $t$  where  $t$  is time. Thus a value of  $t$  would correspond to a point on this curve and you could let  $\mathbf{B}(t)$ ,  $\mathbf{T}(t)$ , and  $\mathbf{N}(t)$  be the binormal, unit tangent, and principal normal at this point of the curve. The following example is typical.

**Example 16.4.5** Given the circular helix,  $\mathbf{R}(t) = (a \cos t)\mathbf{i} + (a \sin t)\mathbf{j} + (bt)\mathbf{k}$ , find the arc length,  $s(t)$ , the unit tangent vector,  $\mathbf{T}(t)$ , the principal normal,  $\mathbf{N}(t)$ , the binormal,  $\mathbf{B}(t)$ , the curvature,  $\kappa(t)$ , and the torsion,  $\tau(t)$ . Here  $t \in [0, T]$ .

The arc length is  $s(t) = \int_0^t (\sqrt{a^2 + b^2}) \, dr = (\sqrt{a^2 + b^2})t$ . Now the tangent vector is obtained using the chain rule as

$$\begin{aligned}\mathbf{T} &= \frac{d\mathbf{R}}{ds} = \frac{d\mathbf{R}}{dt} \frac{dt}{ds} = \frac{1}{\sqrt{a^2 + b^2}} \mathbf{R}'(t) \\ &= \frac{1}{\sqrt{a^2 + b^2}} ((-a \sin t)\mathbf{i} + (a \cos t)\mathbf{j} + b\mathbf{k})\end{aligned}$$

The principal normal:

$$\begin{aligned}\frac{d\mathbf{T}}{ds} &= \frac{d\mathbf{T}}{dt} \frac{dt}{ds} \\ &= \frac{1}{a^2 + b^2} ((-a \cos t)\mathbf{i} + (-a \sin t)\mathbf{j} + 0\mathbf{k})\end{aligned}$$

and so

$$\mathbf{N} = \frac{d\mathbf{T}}{ds} / \left| \frac{d\mathbf{T}}{ds} \right| = -((\cos t)\mathbf{i} + (\sin t)\mathbf{j})$$

The binormal:

$$\begin{aligned}\mathbf{B} &= \frac{1}{\sqrt{a^2 + b^2}} \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ -a \sin t & a \cos t & b \\ -\cos t & -\sin t & 0 \end{vmatrix} \\ &= \frac{1}{\sqrt{a^2 + b^2}} ((b \sin t)\mathbf{i} - b \cos t \mathbf{j} + a\mathbf{k})\end{aligned}$$

Now the curvature,  $\kappa(t) = \left| \frac{d\mathbf{T}}{ds} \right| = \sqrt{\left( \frac{a \cos t}{a^2 + b^2} \right)^2 + \left( \frac{a \sin t}{a^2 + b^2} \right)^2} = \frac{a}{a^2 + b^2}$ . Note the curvature is constant in this example. The final task is to find the torsion. Recall that  $\mathbf{B}' = \tau \mathbf{N}$  where the derivative on  $\mathbf{B}$  is taken with respect to arc length. Therefore, remembering that  $t$  is a function of  $s$ ,

$$\begin{aligned}\mathbf{B}'(s) &= \frac{1}{\sqrt{a^2 + b^2}} ((b \cos t)\mathbf{i} + (b \sin t)\mathbf{j}) \frac{dt}{ds} \\ &= \frac{1}{a^2 + b^2} ((b \cos t)\mathbf{i} + (b \sin t)\mathbf{j}) \\ &= \tau (-(\cos t)\mathbf{i} - (\sin t)\mathbf{j}) = \tau \mathbf{N}\end{aligned}$$

and it follows  $-b/(a^2 + b^2) = \tau$ .

An important application of the usefulness of these ideas involves the decomposition of the acceleration in terms of these vectors of an object moving over a space curve.

**Corollary 16.4.6** Let  $\mathbf{R}(t)$  be a space curve and denote by  $\mathbf{v}(t)$  the velocity,  $\mathbf{v}(t) = \mathbf{R}'(t)$  and let  $v(t) \equiv |\mathbf{v}(t)|$  denote the speed and let  $\mathbf{a}(t)$  denote the acceleration. Then  $\mathbf{v} = v\mathbf{T}$  and  $\mathbf{a} = \frac{dv}{dt}\mathbf{T} + \kappa v^2 \mathbf{N}$ .

**Proof:**  $\mathbf{T} = \frac{d\mathbf{R}}{ds} = \frac{d\mathbf{R}}{dt} \frac{dt}{ds} = \mathbf{v} \frac{dt}{ds}$ . Also,  $s = \int_0^t v(r) \, dr$  and so  $\frac{ds}{dt} = v$  which implies  $\frac{dt}{ds} = \frac{1}{v}$ . Therefore,  $\mathbf{T} = \mathbf{v}/v$  which implies  $\mathbf{v} = v\mathbf{T}$  as claimed.

Now the acceleration is just the derivative of the velocity and so by the Serrat Frenet formulas,

$$\begin{aligned}\mathbf{a} &= \frac{dv}{dt}\mathbf{T} + v\frac{d\mathbf{T}}{dt} \\ &= \frac{dv}{dt}\mathbf{T} + v\frac{d\mathbf{T}}{ds}v = \frac{dv}{dt}\mathbf{T} + v^2\kappa\mathbf{N}\end{aligned}$$

Note how this decomposes the acceleration into a component tangent to the curve and one which is normal to it. Also note that from the above,  $v|\mathbf{T}'| \frac{\mathbf{T}'(t)}{|\mathbf{T}'|} = v^2\kappa\mathbf{N}$  and so  $\frac{|\mathbf{T}'|}{v} = \kappa$  and  $\mathbf{N} = \frac{\mathbf{T}'(t)}{|\mathbf{T}'|}$

From this, it is possible to give an important formula from physics. Suppose an object orbits a point at constant speed,  $v$ . What is the centripetal acceleration of this object? You may know from a physics class that the answer is  $v^2/r$  where  $r$  is the radius. This follows from the above quite easily. The parameterization of the object which is as described is

$$\mathbf{R}(t) = \left( r \cos\left(\frac{v}{r}t\right), r \sin\left(\frac{v}{r}t\right) \right).$$

Therefore,  $\mathbf{T} = \left( -\sin\left(\frac{v}{r}t\right), \cos\left(\frac{v}{r}t\right) \right)$  and  $\mathbf{T}' = \left( -\frac{v}{r}\cos\left(\frac{v}{r}t\right), -\frac{v}{r}\sin\left(\frac{v}{r}t\right) \right)$ . Thus,  $\kappa = |\mathbf{T}'(t)|/v = \frac{1}{r}$ . It follows

$$\mathbf{a} = \frac{dv}{dt}\mathbf{T} + v^2\kappa\mathbf{N} = \frac{v^2}{r}\mathbf{N}.$$

The vector,  $\mathbf{N}$  points from the object toward the center of the circle because it is a positive multiple of the vector,  $\left( -\frac{v}{r}\cos\left(\frac{v}{r}t\right), -\frac{v}{r}\sin\left(\frac{v}{r}t\right) \right)$ .

Sometimes curves don't come to you parametrically. This is unfortunate when it occurs but you can sometimes find a parametric description of such curves.

**Example 16.4.7** Find a parameterization for the intersection of the surfaces  $y + 3z = 2x^2 + 4$  and  $y + 2z = x + 1$ .

You need to solve for  $x$  and  $y$  in terms of  $z$ . This yields

$$z = 2x^2 - x + 3, \quad y = -4x^2 + 3x - 5.$$

Therefore, letting  $t = x$ , the parameterization is  $(x, y, z) = (t, -4t^2 - 5 + 3t, -t + 3 + 2t^2)$ .

**Example 16.4.8** Find a parametrization for the straight line joining  $(3, 2, 4)$  and  $(1, 10, 5)$ .

$(x, y, z) = (3, 2, 4) + t(-2, 8, 1) = (3 - 2t, 2 + 8t, 4 + t)$  where  $t \in [0, 1]$ . Note where this came from. The vector,  $(-2, 8, 1)$  is obtained from  $(1, 10, 5) - (3, 2, 4)$ . Now you should check to see this works.

## 16.5 Exercises

1. Find a parametrization for the intersection of the planes  $2x + y + 3z = -2$  and  $3x - 2y + z = -4$ .
2. Find a parametrization for the intersection of the plane  $3x + y + z = -3$  and the circular cylinder  $x^2 + y^2 = 1$ .
3. Find a parametrization for the intersection of the plane  $4x + 2y + 3z = 2$  and the elliptic cylinder  $x^2 + 4z^2 = 9$ .



4. Find a parametrization for the straight line joining  $(1, 2, 1)$  and  $(-1, 4, 4)$ .
5. Find a parametrization for the intersection of the surfaces  $3y + 3z = 3x^2 + 2$  and  $3y + 2z = 3$ .
6. Let  $\mathbf{R}(t) = (\cos t)\mathbf{i} + (\cos t)\mathbf{j} + (\sqrt{2}\sin t)\mathbf{k}$ . Find the arc length,  $s$  as a function of the parameter,  $t$ , if  $t = 0$  is taken to correspond to  $s = 0$ .
7. Let  $\mathbf{R}(t) = 2\mathbf{i} + (4t + 2)\mathbf{j} + 4t\mathbf{k}$ . Find the arc length,  $s$  as a function of the parameter,  $t$ , if  $t = 0$  is taken to correspond to  $s = 0$ .
8. Let  $\mathbf{R}(t) = e^{5t}\mathbf{i} + e^{-5t}\mathbf{j} + 5\sqrt{2}t\mathbf{k}$ . Find the arc length,  $s$  as a function of the parameter,  $t$ , if  $t = 0$  is taken to correspond to  $s = 0$ .
9. An object moves along the  $x$  axis toward  $(0, 0)$  and then along the curve  $y = x^2$  in the direction of increasing  $x$  at constant speed. Is the force acting on the object a continuous function? Explain.

## 16.6 Independence Of Parameterization

Recall that if  $\mathbf{p}(t) : t \in [a, b]$  was a parameterization of a smooth curve,  $C$ , the length of  $C$  is defined as

$$\int_a^b |\mathbf{p}'(t)| dt$$

If some other parameterization were used to trace out  $C$ , would the same answer be obtained?

**Theorem 16.6.1** *Let  $\phi : [a, b] \rightarrow [c, d]$  be one to one and suppose  $\phi'$  exists and is continuous on  $[a, b]$ . Then if  $f$  is a continuous function defined on  $[a, b]$  which is Riemann integrable<sup>1</sup>,*

$$\int_c^d f(s) ds = \int_a^b f(\phi(t)) |\phi'(t)| dt$$

**Proof:** Let  $F'(s) = f(s)$ . (For example, let  $F(s) = \int_a^s f(r) dr$ .) Then the first integral equals  $F(d) - F(c)$  by the fundamental theorem of calculus. By Lemma 5.7.4 on Page 94,  $\phi$  is either strictly increasing or strictly decreasing. Suppose  $\phi$  is strictly decreasing. Then  $\phi(a) = d$  and  $\phi(b) = c$ . Therefore,  $\phi' \leq 0$  and the second integral equals

$$\begin{aligned} - \int_a^b f(\phi(t)) \phi'(t) dt &= \int_b^a \frac{d}{dt} (F(\phi(t))) dt \\ &= F(\phi(a)) - F(\phi(b)) = F(d) - F(c). \end{aligned}$$

The case when  $\phi$  is increasing is similar. This proves the theorem.

**Lemma 16.6.2** *Let  $\mathbf{f} : [a, b] \rightarrow C$ ,  $\mathbf{g} : [c, d] \rightarrow C$  be parameterizations of a smooth curve which satisfy conditions 1 - 5. Then  $\phi(t) \equiv \mathbf{g}^{-1} \circ \mathbf{f}(t)$  is 1 - 1 on  $(a, b)$ , continuous on  $[a, b]$ , and either strictly increasing or strictly decreasing on  $[a, b]$ .*

**Proof:** It is obvious  $\phi$  is 1 - 1 on  $(a, b)$  from the conditions  $\mathbf{f}$  and  $\mathbf{g}$  satisfy. It only remains to verify continuity on  $[a, b]$  because then the final claim follows from Lemma 5.15.3 on Page 117. If  $\phi$  is not continuous on  $[a, b]$ , then there exists a sequence,  $\{t_n\} \subseteq [a, b]$  such that  $t_n \rightarrow t$  but  $\phi(t_n)$  fails to converge to  $\phi(t)$ . Therefore, for some  $\varepsilon > 0$  there

<sup>1</sup>Recall that all continuous functions of this sort are Riemann integrable.

exists a subsequence, still denoted by  $n$  such that  $|\phi(t_n) - \phi(t)| \geq \varepsilon$ . Using the sequential compactness of  $[c, d]$ , (See Theorem 14.9.5 on Page 347.) there is a further subsequence, still denoted by  $n$  such that  $\{\phi(t_n)\}$  converges to a point,  $s$ , of  $[c, d]$  which is not equal to  $\phi(t)$ . Thus  $\mathbf{g}^{-1} \circ \mathbf{f}(t_n) \rightarrow s$  and still  $t_n \rightarrow t$ . Therefore, the continuity of  $\mathbf{f}$  and  $\mathbf{g}$  imply  $\mathbf{f}(t_n) \rightarrow \mathbf{g}(s)$  and  $\mathbf{f}(t_n) \rightarrow \mathbf{f}(t)$ . Therefore,  $\mathbf{g}(s) = \mathbf{f}(t)$  and so  $s = \mathbf{g}^{-1} \circ \mathbf{f}(t) = \phi(t)$ , a contradiction. Therefore,  $\phi$  is continuous as claimed.

**Theorem 16.6.3** *The length of a smooth curve is not dependent on parameterization.*

**Proof:** Let  $C$  be the curve and suppose  $\mathbf{f} : [a, b] \rightarrow C$  and  $\mathbf{g} : [c, d] \rightarrow C$  both satisfy conditions 1 - 5. Is it true that  $\int_a^b |\mathbf{f}'(t)| dt = \int_c^d |\mathbf{g}'(s)| ds$ ?

Let  $\phi(t) \equiv \mathbf{g}^{-1} \circ \mathbf{f}(t)$  for  $t \in [a, b]$ . Then by the above lemma  $\phi$  is either strictly increasing or strictly decreasing on  $[a, b]$ . Suppose for the sake of simplicity that it is strictly increasing. The decreasing case is handled similarly.

Let  $s_0 \in \phi([a + \delta, b - \delta]) \subset (c, d)$ . Then by assumption 4,  $g'_i(s_0) \neq 0$  for some  $i$ . By continuity of  $g'_i$ , it follows  $g'_i(s) \neq 0$  for all  $s \in I$  where  $I$  is an open interval contained in  $[c, d]$  which contains  $s_0$ . It follows that on this interval,  $g_i$  is either strictly increasing or strictly decreasing. Therefore,  $J \equiv g_i(I)$  is also an open interval and you can define a differentiable function,  $h_i : J \rightarrow I$  by

$$h_i(g_i(s)) = s.$$

As in Theorem 8.1.6 on Page 151, this implies that for  $s \in I$ ,

$$h'_i(g_i(s)) = \frac{1}{g'_i(s)}. \quad (16.4)$$

Now letting  $s = \phi(t)$  for  $s \in I$ , it follows  $t \in J_1$ , an open interval. Also, for  $s$  and  $t$  related this way,  $\mathbf{f}(t) = \mathbf{g}(s)$  and so in particular, for  $s \in I$ ,

$$g_i(s) = f_i(t).$$

Consequently,

$$s = h_i(f_i(t)) = \phi(t)$$

and so, for  $t \in J_1$ ,

$$\phi'(t) = h'_i(f_i(t)) f'_i(t) = h'_i(g_i(s)) f'_i(t) = \frac{f'_i(t)}{g'_i(\phi(t))} \quad (16.5)$$

which shows that  $\phi'$  exists and is continuous on  $J_1$ , an open interval containing  $\phi^{-1}(s_0)$ . Since  $s_0$  is arbitrary, this shows  $\phi'$  exists on  $[a + \delta, b - \delta]$  and is continuous there.

Now  $\mathbf{f}(t) = \mathbf{g} \circ (\mathbf{g}^{-1} \circ \mathbf{f})(t) = \mathbf{g}(\phi(t))$  and it was just shown that  $\phi'$  is a continuous function on  $[a - \delta, b + \delta]$ . It follows

$$\mathbf{f}'(t) = \mathbf{g}'(\phi(t)) \phi'(t)$$

and so, by Theorem 16.6.1,

$$\begin{aligned} \int_{\phi(a+\delta)}^{\phi(b-\delta)} |\mathbf{g}'(s)| ds &= \int_{a+\delta}^{b-\delta} |\mathbf{g}'(\phi(t))| |\phi'(t)| dt \\ &= \int_{a+\delta}^{b-\delta} |\mathbf{f}'(t)| dt. \end{aligned}$$

Now using the continuity of  $\phi$ ,  $\mathbf{g}'$ , and  $\mathbf{f}'$  on  $[a, b]$  and letting  $\delta \rightarrow 0+$  in the above, yields

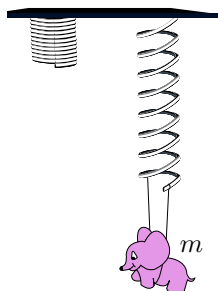
$$\int_c^d |\mathbf{g}'(s)| ds = \int_a^b |\mathbf{f}'(t)| dt$$

and this proves the theorem.

# The Circular Functions Again

The circular functions have already been discussed. You recall the important role of plane geometry in defining them. A major consideration was the careful definition of arc length of a circle and showing the radian measure of an angle is well defined. This involved similar triangles and such notions. On the other hand, the only thing necessary do calculus is the concept of completeness and a way to measure distance, the distance formula. Is it possible to define the circular functions beginning with these things instead of dragging in all that plane geometry? Also, what are some things the circular functions are good for? This chapter is devoted to these questions.

To begin with, there is a physical motivation for the circular functions. Consider a garage door spring. These springs exert a force which resists extension. Attach such a spring to the ceiling and attach a mass,  $m$ , to the bottom end of the spring as shown in the following picture. Any mass will do. It does not have to be a small elephant.



The weight of this mass,  $mg$ , is a downward force which extends the spring, moving the bottom end of the spring downward a distance  $l$  where the upward force exerted by the spring exactly balances the downward force exerted on the mass by gravity. It has been experimentally observed that as long as the extension,  $z$ , of such a spring is not too great, the restoring force exerted by the spring is of the form  $kz$  where  $k$  is some constant which depends on the spring. (It would be different for a slinky than for a garage door spring.) This is known as Hooke's law which is the simplest model for elastic materials. Therefore,  $mg = kl$ . Now let  $y$  be the displacement from this equilibrium position of the bottom of the spring with the positive direction being up. Thus the acceleration of the spring would be  $y''$ . The extension of the spring in terms of  $y$  would be  $(l - y)$ . Then Newton's second law along with Hooke's law imply

$$my'' = k(l - y) - mg$$

and since  $kl - mg = 0$ , this yields

$$my'' + ky = 0.$$

Dividing by  $m$  and letting  $\omega^2 = k/m$  yields the equation for undamped oscillation,

$$y'' + \omega^2 y = 0.$$

Based on physical reasoning just presented, there should be a solution to this equation. It is the displacement of the bottom end of a spring from the equilibrium position. However, it is not enough to base questions of existence in mathematics on physical intuition, although it is sometimes done. The following theorem gives the necessary existence and uniqueness results.

**Theorem 17.0.4** *There exists exactly one function,  $y$ , which satisfies the following initial value problem.*

$$y'' + \omega^2 y = 0, \quad y(a) = y_0, y'(a) = y_1. \quad (17.1)$$

**Proof:** First consider the question of uniqueness. Suppose both  $z$  and  $y$  are solutions to the problem (17.1). Then let  $x = z - y$  and use Theorem 6.2.6 to conclude  $x$  satisfies the following initial value problem.

$$x'' + \omega^2 x = 0, \quad x(a) = x'(a) = 0.$$

Now multiply both sides of the differential equation by  $x'$  and use Theorem 6.2.6 again to write

$$\frac{d}{dt} \left( \frac{(x')^2}{2} + \frac{\omega^2 (x)^2}{2} \right) = 0.$$

By Corollary 6.8.4 on Page 133 the expression

$$\frac{(x')^2}{2} + \frac{\omega^2 (x)^2}{2}$$

must equal a constant. However this expression equals zero when  $t = a$  and so it must be zero for all  $t$ . Therefore, in particular,  $x(t) = 0$  and so  $z(t) = y(t)$ . This proves the uniqueness part.

To establish the existence of a solution, let  $z(t) \equiv y(t + a)$ . Then  $y$  solves (17.1) if and only if  $z$  solves

$$z'' + \omega^2 z = 0, \quad z(0) = y_0, z'(0) = y_1. \quad (17.2)$$

Now let

$$C(t) \equiv \sum_{n=0}^{\infty} (-1)^n \frac{t^{2n}}{(2n)!}, \quad S(t) \equiv \sum_{n=0}^{\infty} (-1)^n \frac{t^{2n+1}}{(2n+1)!}. \quad (17.3)$$

Using the ratio test, you see both series in the above converge for all  $t \in \mathbb{R}$ . (Of course you know what these series are already but you don't need to know this for the purposes of this presentation. Remember, you can define functions in terms of power series.) Then you see right away using the theorem about differentiation of power series that

$$S'(t) = C(t) \quad \text{and} \quad C'(t) = -S(t). \quad (17.4)$$

Consider the function,  $z(t) = y_0 C(\omega t) + \frac{y_1}{\omega} S(\omega t)$ . It follows from the chain rule and what was just observed about  $S'$  and  $C'$  that  $z$  solves (17.2). Therefore, our solution to (17.1) is the function,  $y(t) \equiv z(t - a)$ .

The quantity which is a constant in the above proof of uniqueness is sometimes referred to as the energy.

It follows immediately from (17.3) that  $C(t)$  satisfies the following initial value problems.

$$C'' + C = 0, \quad C(0) = 1, C'(0) = 0.$$

In terms of the spring, this would have the spring constant,  $k$  equal in magnitude to the mass,  $m$  and the function would be the displacements from equilibrium of the bottom of the spring when initially it is raised one unit above the equilibrium point and released with zero velocity. The second of these functions,  $S(t)$  is the solution to the initial value problem

$$S'' + S = 0, S(0) = 0, S'(0) = 1.$$

In terms of the spring, the spring constant,  $k$  would be equal in magnitude to the mass,  $m$  and the function would be the displacements from equilibrium of the bottom of the spring when initially it is at the equilibrium point and moving upward with unit speed. The next theorem is on properties of these functions,  $C$  and  $S$ . These properties will be obtained directly from the initial value problems they satisfy. In addition, another proof is given of Property (17.4) which is based on the initial value problems satisfied by the functions rather than the differentiation theorem for power series.

**Theorem 17.0.5** *The functions  $C$  and  $S$  satisfy the following properties.*

$$C^2(t) + S^2(t) = 1, \quad (17.5)$$

$$S(-t) = -S(t), C(-t) = C(t), \quad (17.6)$$

*These conditions in (17.6) say that  $S$  is odd and  $C$  is even because the property of being odd or even is defined in terms of these equations.*

$$S(x+y) = S(x)C(y) + C(x)S(y), \quad (17.7)$$

$$C(x-y) = C(x)C(y) + S(x)S(y), \quad (17.8)$$

For  $x \geq 0$ ,

$$x - \frac{x^3}{6} \leq S(x), \quad (17.9)$$

$$1 - \frac{x^2}{2} \leq C(x) \leq 1 - \frac{x^2}{2} + \frac{x^4}{24}. \quad (17.10)$$

**Proof:** To verify (17.4) consider the function,  $y(x) = C' + S$ . Then

$$y'' + y = C''' + C' + S'' + S = 0$$

and

$$y(0) = C'(0) + S(0) = 0 + 0 = 0,$$

while

$$y'(0) = C''(0) + S'(0) = -C(0) + C(0) = 0.$$

By the Theorem 17.0.4,  $y(x) = 0$ . The second half of (17.4) is proved similarly only this time consider  $y(x) = S'(x) - C(x)$ .

It will be shown that  $y$  solves the differential equation and the initial condition,  $y(0) = S'(0) - C(0) = 1 - 1 = 0$ . Thus by Theorem 17.0.4 again  $y(x) = 0$ .

To verify (17.5) use the initial value problem satisfied by  $C$  to write

$$C'' + C = 0.$$

Now multiply by  $C'$  and use Theorem 6.2.6 to obtain

$$\frac{d}{dt} \left( \frac{(C'(t))^2}{2} + \frac{(C(t))^2}{2} \right) = 0$$

and so

$$\frac{(C'(t))^2}{2} + \frac{(C(t))^2}{2} = c$$

for some constant,  $c$ . Now the initial conditions for  $C$  imply this constant,  $c$  equals  $1/2$ . Using  $C' = -S$ , yields (17.5).

Let  $y(t) = S(-t) + S(t)$ . Then by Theorem 6.2.6  $y'(t) = -S'(-t) + S'(t)$  and so  $y(0) = 0$  and  $y'(0) = -1 + 1 = 0$ . Also using Theorem 6.2.6 again,  $y''(t) = S'(-t) + S'(t)$  and so  $y'' + y = 0$ . Then by Theorem 17.0.4, it follows  $y(t) = 0$ . This proves the first part of (17.6). The second part is similar.

To verify (17.7), fix  $y$  and let  $y(x) = S(x+y) - (S(x)C(y) + S(y)C(x))$ . Then from (17.4),

$$y'(x) = C(x+y) - (C(x)C(y) - S(y)S(x))$$

and

$$y''(x) = -S(x+y) - (-S(x)C(y) - S(y)C(x))$$

so

$$y'' + y = 0.$$

Also  $y(0) = S(y) - S(y) = 0$  and  $y'(0) = C(y) - C(y) = 0$ . Therefore, by Theorem 17.0.4, it follows  $y(x) = 0$  proving (17.7). The proof of (17.8) is completely similar and is left as an exercise.

This brings us to the estimates. Let  $w(x) = x - S(x)$ . Then  $w(0) = 0$  and

$$w'(x) = 1 - C(x) \geq 0.$$

By Corollary 6.8.5 on Page 133, it follows that  $w$  is increasing on  $[0, \infty)$ . Therefore,  $w(x) \geq 0$  which shows that

$$x - S(x) \geq 0.$$

Let  $z(x) = C(x) - \left(1 - \frac{x^2}{2}\right)$ . Then  $z(0) = 0$  and

$$z'(x) = -S(x) + x \geq 0$$

so by Corollary 6.8.5 on Page 133, it follows that  $z(x) = C(x) - \left(1 - \frac{x^2}{2}\right) \geq 0$  for all  $x \geq 0$ , proving the bottom half of (17.10). Now let  $y(x) = S(x) - \left(x - \frac{x^3}{6}\right)$ . Then  $y(0) = 0$  and

$$y'(x) = C(x) - \left(1 - \frac{x^2}{2}\right) \geq 0$$

so by Corollary 6.8.5, it follows that  $y(x) = S(x) - \left(x - \frac{x^3}{6}\right) \geq 0$  for all  $x \geq 0$ . This proves (17.9). Now let  $u(x) = \left(1 - \frac{x^2}{2} + \frac{x^4}{24}\right) - C(x)$ . Then  $u(0) = 0$  and

$$u'(x) = -x + \frac{x^3}{6} + S(x) \geq 0$$

by (17.9). Therefore, by Corollary 6.8.5,  $u(x) \geq 0$  for all  $x \geq 0$  and this shows the top half of (17.10).

With (17.9) and (17.10) other properties of  $C$  and  $S$  related to their periodic behavior may be obtained.

**Theorem 17.0.6** *There exists a positive number,  $P$  such that  $C(x) > 0$  on  $[0, \frac{P}{2})$  and  $C(\frac{P}{2}) = 0$ . In addition to this, both  $C$  and  $S$  are periodic of period  $2P$ . This means that for all  $x \in \mathbb{R}$ ,*

$$C(x + 2P) = C(x), \quad S(x + 2P) = S(x).$$

*Also the function,  $\mathbf{r}(t) \equiv (C(t), S(t))$  is one to one for  $t \in [0, 2P)$  and  $C(t) = \cos t$  while  $S(t) = \sin t$  for  $t$  given in radians and  $P = \pi$  where  $2\pi$  is the circumference of the unit circle.*

**Proof:** From (17.10) that  $C(2) < 0$  because  $C(2) \leq \left(1 - \frac{2^2}{2} + \frac{2^4}{24}\right) = -\frac{1}{3} < 0$ . Recall,  $C(1) = 1 > 0$ . By the intermediate value theorem there exists a number,  $P/2 \in (0, 2)$  such that  $C(P/2) = 0$ . There is only one such number because for  $x \in (0, 2)$ ,

$$C'(x) = -S'(x) \leq -x + \frac{x^3}{6} = (-x) \left(1 - \frac{x^2}{6}\right) < 0$$

showing that  $C$  is strictly decreasing on this interval because of Corollary 6.8.6. This proves the first part of the theorem.  $S$  is increasing on  $(0, P/2)$  because  $S'(x) = C(x)$  and  $C(x) > 0$  on  $(0, P/2)$ . Therefore, from (17.5),  $S(P/2) = 1$ . Also this shows that  $\mathbf{r}(t)$  is one to one for  $t \in [0, P/2]$ . Now consider what happens on the interval,  $[P/2, P]$ . Points on this interval are of the form,  $x + P/2$  where  $x \in [0, P/2]$ . From (17.6), (17.7), and (17.8),

$$\begin{aligned} C(x + P/2) &= C(x)C(P/2) - S(x)S(P/2) \\ &= -S(x) \leq 0, < 0 \text{ if } x \in (0, P/2) \end{aligned}$$

while

$$\begin{aligned} S(x + P/2) &= S(x)C(P/2) + S(P/2)C(x) \\ &= C(x) \geq 0, > 0 \text{ if } x \in (0, P/2). \end{aligned}$$

Therefore,  $S$  is strictly decreasing on  $[P/2, P]$  because its derivative,  $C$  is negative on  $(P/2, P)$ . Also  $C$  is strictly decreasing on this interval because its derivative,  $-S$  is negative on  $(P/2, P)$ . Since  $C$  is strictly decreasing on the interval  $[0, P]$  it follows that  $\mathbf{r}$  is one to one on  $[0, P]$ . Furthermore,

$$\begin{aligned} S(P) &= S(P/2)C(P/2) + S(P/2)C(P/2) = 0 \\ C(P) &= C(P/2)C(P/2) - S(P/2)S(P/2) = -1. \end{aligned}$$

In going from  $P$  to  $3P/2$ , use similar arguments to verify that  $S$  is strictly decreasing and that  $C$  is strictly increasing and that  $S(3P/2) = -1$  while  $C(-3P/2) = 0$ . If  $\mathbf{r}(t_1) = \mathbf{r}(t_2)$  for  $t_1 \in [0, P]$  and  $t_2 \in [P, 3P/2]$ , then  $C(t_1) = C(t_2)$  and  $S(t_1) = S(t_2)$ . Therefore,  $S(t_2) \leq 0$  and so  $S(t_1) \leq 0$ . But  $S(t_1) \geq 0$  and so  $t_1 = P = t_2$ . It follows that  $\mathbf{r}$  is one to one on  $[0, 3P/2]$ . In going from  $3P/2$  to  $2P$ , use the same arguments to verify that  $C$  and  $S$  are strictly increasing and  $C(2P) = 1$  while  $S(2P) = 0$ . Now if  $\mathbf{r}(t_1) = \mathbf{r}(t_2)$  for  $t_2 \in (3P/2, 2P)$  and  $t_1 \in [0, 3P/2]$ ,  $C(t_1) = C(t_2) > 0$  and so,  $t_1 \in [0, P/2)$ . Now  $S(t_1) = S(t_2) < 0$  which can't happen for  $t_1 \in [0, P/2]$ . Therefore,  $\mathbf{r}$  is one to one on  $[0, 2P)$ . It only remains to verify that  $C$  and  $S$  are periodic of period  $2P$ , that  $P = \pi$ , and that  $C = \cos$  while  $S = \sin$ .

First consider the claim that  $S$  and  $C$  are  $2P$  periodic.

$$\begin{aligned} S(x + 2P) &= S(x + P)C(P) + C(x + P)S(P) \\ &= -S(x + P) \\ &= -[S(x)C(P) + C(x)S(P)] = S(x). \end{aligned}$$

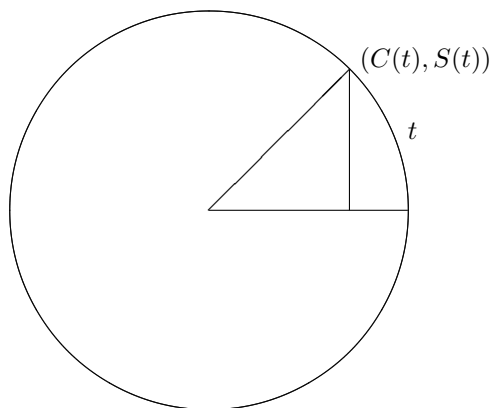
To verify the assertion about  $C$ ,

$$\begin{aligned} C(x+2P) &= C(x+P)C(P) - S(x+P)S(P) \\ &= -C(x+P) \\ &= -[C(x)C(P) - S(x)S(P)] = C(x). \end{aligned}$$

Now it remains to verify the other assertions. First note that by (17.5),  $\mathbf{r}(t) = (C(t), S(t))$  is a point on the unit circle. Since  $\mathbf{r}$  is one to one on  $(0, 2P)$ , the length of the path traced out by  $\mathbf{r}$  for  $t \in [0, s]$  is given by  $A(s)$  where  $A$  solves the initial value problem,

$$A'(t) = \sqrt{(C'(t))^2 + (S'(t))^2} = \sqrt{S^2(t) + C^2(t)} = 1, \quad A(0) = 0.$$

The solution to this is just  $A(t) = t$  and so the length of this path traced out by  $\mathbf{r}$  is just  $s$ . In other words,  $(C(t), S(t))$  is the point on the unit circle which is at a distance of  $t$  measured along an arc which starts at  $(1, 0)$  and proceeds counter clock wise to this point as shown in the following picture in which  $t$  denotes the length of the indicated arc of the circle.



Thus, drawing a right triangle with the radius of the circle as its hypotenuse and  $t$  the angle just described,  $C(t) = \cos t$  and  $S(t) = \sin t$ . Now the length of the curve traced out by  $\mathbf{r}$  for  $t \in [0, 2P]$  is just  $2P$  but from the above discussion, this involves starting at  $(1, 0)$  and continuing around the unit circle till  $(1, 0)$ . Define  $2\pi$  to be this length. Therefore,  $P = \pi$ . This proves the theorem.

Note that this finishes off all the trig identities as well thanks to the previous theorem, so this gives an alternate way of defining the circular functions. No reference was made to plane geometry, just the definition of distance along with the notion of arc length.

## 17.1 The Equations Of Undamped And Damped Oscillation

With the circular functions, the equation of undamped oscillation,  $y'' + \omega^2 y = 0$  may be solved. This equation resulted from the oscillations of a heavy weight on a spring but it occurs in many other physical contexts also.

**Theorem 17.1.1** *The initial value problem,*

$$y'' + \omega^2 y = 0, \quad y(0) = y_0, y'(0) = y_1 \quad (17.11)$$



has a unique solution and this solution is

$$y(t) = y_0 \cos(\omega t) + \frac{y_1}{\omega} \sin(\omega t). \quad (17.12)$$

**Proof:** You should verify that (17.12) does indeed provide a solution to the initial value problem (17.11). It follows from Theorem 17.0.4 that this is the only solution to this problem.

Now consider another sort of differential equation,

$$y'' - a^2 y = 0, \quad a > 0 \quad (17.13)$$

To give the complete solution, let  $Dy \equiv y'$ . Then the differential equation may be written as

$$(D + a)(D - a)y = 0.$$

Let  $z = (D - a)y$ . Thus  $(D + a)z = 0$  and so  $z(t) = C_1 e^{-at}$  from Theorem 9.19.1 on Page 221. Therefore,

$$(D - a)y \equiv y' - ay = C_1 e^{-at}.$$

Multiply both sides of this last equation by  $e^{-at}$ . By the product and chain rules,

$$\frac{d}{dt}(e^{-at}y) = C_1 e^{-2at}.$$

Therefore,

$$e^{-at}y = \frac{C_1}{-2a} e^{-2at} + C_2$$

and so

$$y = \frac{C_1}{-2a} e^{-at} + C_2 e^{at}.$$

Now since  $C_1$  is arbitrary, it follows any solution of (17.13) is of the form  $y = C_1 e^{-at} + C_2 e^{at}$ . Now you should verify that any expression of this form actually solves the equation, (17.13). This proves most of the following theorem.

**Theorem 17.1.2** *Every solution of the differential equation,  $y'' - a^2 y = 0$  is of the form  $C_1 e^{-at} + C_2 e^{at}$  for some constants  $C_1$  and  $C_2$  provided  $a > 0$ . In the case when  $a = 0$ , every solution of  $y'' = 0$  is of the form  $C_1 t + C_2$  for some constants,  $C_1$  and  $C_2$ .*

All that remains of the proof is to do the part when  $a = 0$  which is left as an exercise involving the mean value theorem.

Now consider the differential equation of damped oscillation. In the example of the object bobbing on the end of a spring,

$$my'' = -ky$$

where  $k$  was the spring constant and  $m$  was the mass of the object. Suppose the object is also attached to a dash pot. This is a device which resists motion like a shock absorber on a car. You know how these work. If the car is just sitting still the shock absorber applies no force to the car. It only gives a force in response to up and down motion of the car and you assume this force is proportional to the velocity and opposite the velocity. Thus in our spring example, you would have

$$my'' = -ky - \delta^2 y'$$

where  $\delta^2$  is the constant of proportionality of the resisting force. Dividing by  $m$  and adjusting the coefficients, such damped oscillation satisfies an equation of the form,

$$y'' + 2by' + ay = 0. \quad (17.14)$$

Actually this is a general homogeneous second order equation, more general than what results from damped oscillation. Concerning the solutions to this equation, the following theorem is given. In this theorem the first case is referred to as the underdamped case. The second case is called the critically damped case and the third is called the overdamped case.

**Theorem 17.1.3** Suppose  $b^2 - a < 0$ . Then all solutions of (17.14) are of the form

$$e^{-bt} (C_1 \cos(\omega t) + C_2 \sin(\omega t)). \quad (17.15)$$

where  $\omega = \sqrt{a - b^2}$  and  $C_1$  and  $C_2$  are constants. In the case that  $b^2 - a = 0$  the solutions of (17.14) are of the form

$$e^{-bt} (C_1 + C_2 t). \quad (17.16)$$

In the case that  $b^2 - a > 0$  the solutions are of the form

$$e^{-bt} (C_1 e^{-rt} + C_2 e^{rt}), \quad (17.17)$$

where  $r = \sqrt{b^2 - a}$ .

**Proof:** Let  $z = e^{bt}y$  and write (17.14) in terms of  $z$ . Thus,  $z$  is a solution to the equation

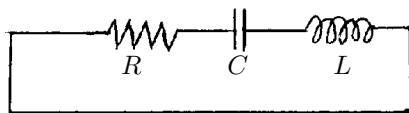
$$z'' + (a - b^2)z = 0. \quad (17.18)$$

If  $b^2 - a < 0$ , then by Theorem 17.1.1,  $z(t) = C_1 \cos(\omega t) + C_2 \sin(\omega t)$  where  $\omega = \sqrt{a - b^2}$ . Therefore,

$$y = e^{-bt} (C_1 \cos(\omega t) + C_2 \sin(\omega t))$$

as claimed. The other two cases are completely similar. They use Theorem 17.1.2 rather than Theorem 17.1.1.

**Example 17.1.4** An important example of these equations occurs in an electrical circuit having a capacitor, a resistor, and an inductor in series as shown in the following picture.



The voltage drop across the inductor is  $L \frac{di}{dt}$  where  $i$  is the current and  $L$  is the inductance. The voltage drop across the resistor is  $Ri$  where  $R$  is the resistance. This is according to Ohm's law. The voltage drop across the capacitor is  $v = \frac{Q}{C}$  where  $Q$  is the charge on the capacitor and  $C$  is a constant called the capacitance. The current equals the rate of change of the charge on the capacitor. Thus  $i = Q' = Cv'$ . When these voltages are summed, you must get zero because there is no voltage source in the circuit. Thus  $L \frac{di}{dt} + Ri + \frac{Q}{C} = 0$  and written in terms of the voltage drop across the capacitor, this becomes  $LCv'' + CRv' + v = 0$ , a second order linear differential equation of the sort discussed above.

## 17.2 Exercises

1. Prove from the initial value problems satisfied by  $S$  and  $C$  the formula (17.8). Also verify the second half of (17.6).
2. Verify that  $y = C_1 e^{-at} + C_2 e^{at}$  solves the differential equation, (17.13).
3. Verify (17.18).
4. Verify that all solutions to the differential equation,  $y'' = 0$  are of the form  $y = C_1 t + C_2$ .
5. Show from Corollary 6.8.5 and Theorem 17.0.5 that for all  $x \geq 0$ ,

$$S(x) \leq x - \frac{x^3}{6} + \frac{x^5}{120}. \quad (17.19)$$

6. Using Theorem 17.0.5 and Problem 5, estimate  $S(.1)$  and  $C(.1)$ . Give upper and lower bounds for these numbers.
7. Using Theorem 17.0.5 and Problem 5 and Theorem 5.9.5, establish the limit,

$$\lim_{x \rightarrow 0} \frac{S(x)}{x} = 1.$$

8. A mass of ten Kilograms is suspended from a spring attached to the ceiling. This mass causes the end of the spring to be displaced a distance of 39.2 *cm*. The mass end of the spring is then pulled down a distance of one *cm*. and released. Find the displacement from the equilibrium position of the end of the spring as a function of time. Assume the acceleration of gravity is  $9.8m/\text{sec}^2$ .
9. Keep everything the same in Problem 8 except suppose the suspended end of the spring is also attached to a dash pot which provides a force opposite the direction of the velocity having magnitude  $10\sqrt{19}|v|$  Newtons for  $|v|$  the speed. Give the displacement as before.
10. In (17.14) consider the equation in which  $b = -1$  and  $a = 3$ . Explain why this equation describes a physical system which has some dubious properties.
11. Solve the initial value problem,  $y'' + 5y' - y = 0$ ,  $y(0) = 1$ ,  $y'(0) = 0$ .
12. Solve the initial value problem  $y'' + 2y' + 2y = 0$ ,  $y(0) = 0$ ,  $y'(0) = 1$ .
13. In the case of undamped oscillation show the solution can be written in the form  $A \cos(\omega t - \phi)$  where  $\phi$  is some angle called a phase shift and a constant,  $A$ , called the amplitude.
14. Using Theorem 17.0.5 and Problem 5 and Theorem 5.9.5, establish the limit,

$$\lim_{x \rightarrow 0} \frac{1 - C(x)}{x} = 0.$$

15. Is the derivative of a function always continuous? **Hint:** Consider

$$h(x) = \begin{cases} x^2 \sin\left(\frac{1}{x}\right) & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases}.$$

Show  $h'(0) = 0$ . What is  $h'(x)$  for  $x \neq 0$ ?

16. Suppose  $f(t) = a \cos \omega t + b \sin \omega t$ . Show there exists  $\phi$  such that

$$f(t) = \sqrt{a^2 + b^2} \sin(\omega t + \phi).$$

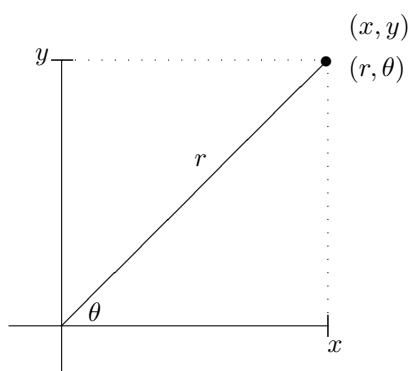
**Hint:**  $f(t) = \sqrt{a^2 + b^2} \left( \frac{a}{\sqrt{a^2 + b^2}} \cos \omega t + \frac{b}{\sqrt{a^2 + b^2}} \sin \omega t \right)$  and  $\left( \frac{b}{\sqrt{a^2 + b^2}}, \frac{a}{\sqrt{a^2 + b^2}} \right)$  is a point on the unit circle so it is of the form  $(\cos \phi, \sin \phi)$ . Now recall the formulas for the sine and cosine of sums of angles.

17. For  $\theta = 0, \pi/6, \pi/4, \pi/3, \pi/2, 2\pi/3, 3\pi/4, 5\pi/6, \pi, 7\pi/6, 5\pi/4, 4\pi/3, 3\pi/2, 5\pi/3, 7\pi/4, 11\pi/6, 2\pi$ , show these angles on the unit circle and find the sine and cosine of each one. These are important examples. Do this without any reference to triangles or plane geometry.

# Curvilinear Coordinate Systems

## 18.1 Polar Cylindrical And Spherical Coordinates

So far points have been identified in terms of Cartesian coordinates but there are other ways of specifying points in two and three dimensional space. These other ways involve using a list of two or three numbers which have a totally different meaning than Cartesian coordinates to specify a point in two or three dimensional space. In general these lists of numbers which have a different meaning than Cartesian coordinates are called Curvilinear coordinates. Probably the simplest curvilinear coordinate system is that of polar coordinates. The idea is suggested in the following picture.



You see in this picture, the number  $r$  identifies the distance of the point from the origin,  $(0,0)$  while  $\theta$  is the angle shown between the positive  $x$  axis and the line from the origin to the point. This angle will always be given in radians and is in the interval  $[0, 2\pi)$ . Thus the given point, indicated by a small dot in the picture, can be described in terms of the Cartesian coordinates,  $(x, y)$  or the polar coordinates,  $(r, \theta)$ . How are the two coordinates systems related? From the picture,

$$x = r \cos(\theta), \quad y = r \sin(\theta). \quad (18.1)$$

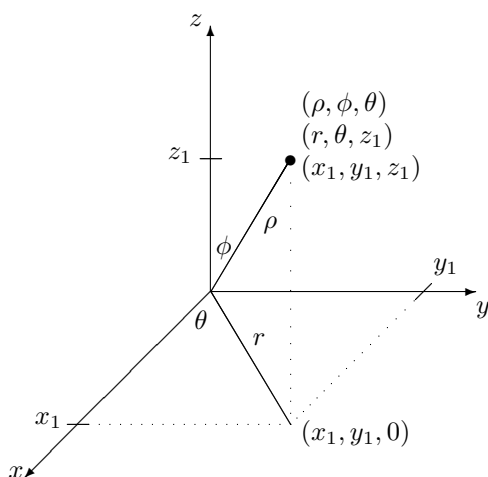
**Example 18.1.1** *The polar coordinates of a point in the plane are  $(5, \frac{\pi}{6})$ . Find the Cartesian or rectangular coordinates of this point.*

From (18.1),  $x = 5 \cos(\frac{\pi}{6}) = \frac{5}{2}\sqrt{3}$  and  $y = 5 \sin(\frac{\pi}{6}) = \frac{5}{2}$ . Thus the Cartesian coordinates are  $(\frac{5}{2}\sqrt{3}, \frac{5}{2})$ .

**Example 18.1.2** *Suppose the Cartesian coordinates of a point are  $(3, 4)$ . Find the polar coordinates.*

Recall that  $r$  is the distance from  $(0, 0)$  and so  $r = 5 = \sqrt{3^2 + 4^2}$ . It remains to identify the angle. Note the point is in the first quadrant, (Both the  $x$  and  $y$  values are positive.) Therefore, the angle is something between  $0$  and  $\pi/2$  and also  $3 = 5 \cos(\theta)$ , and  $4 = 5 \sin(\theta)$ . Therefore, dividing yields  $\tan(\theta) = 4/3$ . At this point, use a calculator or a table of trigonometric functions to find that at least approximately,  $\theta = .927295$  radians.

Now consider two three dimensional generalizations of polar coordinates. The following picture serves as motivation for the definition of these two other coordinate systems.



In this picture,  $\rho$  is the distance between the origin, the point whose Cartesian coordinates are  $(0, 0, 0)$  and the point indicated by a dot and labeled as  $(x_1, y_1, z_1)$ ,  $(r, \theta, z_1)$ , and  $(\rho, \phi, \theta)$ . The angle between the positive  $z$  axis and the line between the origin and the point indicated by a dot is denoted by  $\phi$ , and  $\theta$ , is the angle between the positive  $x$  axis and the line joining the origin to the point  $(x_1, y_1, 0)$  as shown, while  $r$  is the length of this line. Thus  $r$  and  $\theta$  determine a point in the plane determined by letting  $z = 0$  and  $r$  and  $\theta$  are the usual polar coordinates. Thus  $r \geq 0$  and  $\theta \in [0, 2\pi)$ . Letting  $z_1$  denote the usual  $z$  coordinate of a point in three dimensions, like the one shown as a dot,  $(r, \theta, z_1)$  are the cylindrical coordinates of the dotted point. The spherical coordinates are determined by  $(\rho, \phi, \theta)$ . When  $\rho$  is specified, this indicates that the point of interest is on some sphere of radius  $\rho$  which is centered at the origin. Then when  $\phi$  is given, the location of the point is narrowed down to a circle and finally,  $\theta$  determines which point is on this circle. Let  $\phi \in [0, \pi]$ ,  $\theta \in [0, 2\pi)$ , and  $\rho \in [0, \infty)$ . The picture shows how to relate these new coordinate systems to Cartesian coordinates. For Cylindrical coordinates,

$$\begin{aligned}x &= r \cos(\theta), \\y &= r \sin(\theta), \\z &= z\end{aligned}$$

and for spherical coordinates,

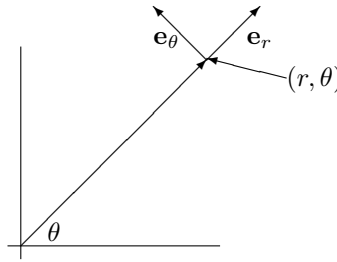
$$\begin{aligned}x &= \rho \sin(\phi) \cos(\theta), \\y &= \rho \sin(\phi) \sin(\theta), \\z &= \rho \cos(\phi).\end{aligned}$$

Spherical coordinates should be especially interesting to you because you live on the surface of a sphere. This has been known for several hundred years. You may also know

that the standard way to determine position on the earth is to give the longitude and latitude. The latitude corresponds to  $\phi$  and the longitude corresponds to  $\theta$ .<sup>1</sup>

## 18.2 The Acceleration In Polar Coordinates

Sometimes you have information about forces which act not in the direction of the coordinate axes but in some other direction. When this is the case, it is often useful to express things in terms of different coordinates which are consistent with these directions. A good example of this is the force exerted by the sun on a planet. This force is always directed toward the sun and so the force vector changes as the planet moves. To discuss this, consider the following simple diagram in which two unit vectors,  $\mathbf{e}_r$  and  $\mathbf{e}_\theta$  are shown.



The vector,  $\mathbf{e}_r = (\cos \theta, \sin \theta)$  and the vector,  $\mathbf{e}_\theta = (-\sin \theta, \cos \theta)$ . You should convince yourself that the picture above corresponds to this definition of the two vectors. Note that  $\mathbf{e}_r$  is a unit vector pointing away from  $\mathbf{0}$  and

$$\mathbf{e}_\theta = \frac{d\mathbf{e}_r}{d\theta}, \quad \mathbf{e}_r = -\frac{d\mathbf{e}_\theta}{d\theta}. \quad (18.2)$$

Now consider the position vector from  $\mathbf{0}$  of a point in the plane,  $\mathbf{r}(t)$ . Then

$$\mathbf{r}(t) = r(t) \mathbf{e}_r(\theta(t))$$

where  $r(t) = |\mathbf{r}(t)|$ . Thus  $r(t)$  is just the distance from the origin,  $\mathbf{0}$  to the point. What is the velocity and acceleration? Using the chain rule,

$$\frac{d\mathbf{e}_r}{dt} = \frac{d\mathbf{e}_r}{d\theta} \theta'(t), \quad \frac{d\mathbf{e}_\theta}{dt} = \frac{d\mathbf{e}_\theta}{d\theta} \theta'(t)$$

and so from (18.2),

$$\frac{d\mathbf{e}_r}{dt} = \theta'(t) \mathbf{e}_\theta, \quad \frac{d\mathbf{e}_\theta}{dt} = -\theta'(t) \mathbf{e}_r. \quad (18.3)$$

Using (18.3) as needed along with the product rule and the chain rule,

$$\begin{aligned} \mathbf{r}'(t) &= r'(t) \mathbf{e}_r + r(t) \frac{d}{dt} (\mathbf{e}_r(\theta(t))) \\ &= r'(t) \mathbf{e}_r + r(t) \theta'(t) \mathbf{e}_\theta. \end{aligned}$$

Next consider the acceleration.

$$\begin{aligned} \mathbf{r}''(t) &= r''(t) \mathbf{e}_r + r'(t) \frac{d\mathbf{e}_r}{dt} + r'(t) \theta'(t) \mathbf{e}_\theta + r(t) \theta''(t) \mathbf{e}_\theta + r(t) \theta'(t) \frac{d}{dt} (\mathbf{e}_\theta) \\ &= r''(t) \mathbf{e}_r + 2r'(t) \theta'(t) \mathbf{e}_\theta + r(t) \theta''(t) \mathbf{e}_\theta + r(t) \theta'(t) (-\mathbf{e}_r) \theta'(t) \\ &= \left( r''(t) - r(t) \theta'(t)^2 \right) \mathbf{e}_r + \left( 2r'(t) \theta'(t) + r(t) \theta''(t) \right) \mathbf{e}_\theta. \end{aligned} \quad (18.4)$$

<sup>1</sup>Actually latitude is determined on maps and in navigation by measuring the angle from the equator rather than the pole but it is essentially the same idea that we have presented here.

This is a very profound formula. Consider the following examples.

**Example 18.2.1** Suppose an object of mass  $m$  moves at a uniform speed,  $s$ , around a circle of radius  $R$ . Find the force acting on the object.

By Newton's second law, the force acting on the object is  $m\mathbf{r}''$ . In this case,  $r(t) = R$ , a constant and since the speed is constant,  $\theta'' = 0$ . Therefore, the term in (18.4) corresponding to  $\mathbf{e}_\theta$  equals zero and  $m\mathbf{r}'' = -R\theta'(t)^2 \mathbf{e}_r$ . The speed of the object is  $s$  and so it moves  $s/R$  radians in unit time. Thus  $\theta'(t) = s/R$  and so

$$m\mathbf{r}'' = -mR \left( \frac{s}{R} \right)^2 \mathbf{e}_r = -m \frac{s^2}{R} \mathbf{e}_r.$$

This is the familiar formula for centripetal force from elementary physics, obtained as a very special case of (18.4).

**Example 18.2.2** A platform rotates at a constant speed in the counter clockwise direction and an object of mass  $m$  moves from the center of the platform toward the edge at constant speed. What forces act on this object?

Let  $v$  denote the constant speed of the object moving toward the edge of the platform. Then

$$r'(t) = v, r''(t) = 0, \theta''(t) = 0,$$

while  $\theta'(t) = \omega$ , a positive constant. From (18.4)

$$m\mathbf{r}''(t) = -mr(t)\omega^2 \mathbf{e}_r + m2v\omega \mathbf{e}_\theta.$$

Thus the object experiences centripetal force from the first term and also a funny force from the second term which is in the direction of rotation of the platform. You can observe this by experiment if you like. Go to a playground and have someone spin one of those merry go rounds while you ride it and move from the center toward the edge. The term  $2r'\theta'$  is called the Coriolis force.

Suppose at each point of space,  $\mathbf{r}$  is associated a force,  $\mathbf{F}(\mathbf{r})$  which a given object of mass  $m$  will experience if its position vector is  $\mathbf{r}$ . This is called a force field. a force field is a central force field if  $\mathbf{F}(\mathbf{r}) = g(\mathbf{r}) \mathbf{e}_r$ . Thus in a central force field, the force an object experiences will always be directed toward or away from the origin,  $\mathbf{0}$ . The following simple lemma is very interesting because it says that in a central force field, objects must move in a plane.

**Lemma 18.2.3** Suppose an object moves in three dimensions in such a way that the only force acting on the object is a central force. Then the motion of the object is in a plane.

**Proof:** Let  $\mathbf{r}(t)$  denote the position vector of the object. Then from the definition of a central force and Newton's second law,

$$m\mathbf{r}'' = g(\mathbf{r}) \mathbf{r}.$$

Therefore,  $m\mathbf{r}'' \times \mathbf{r} = m(\mathbf{r}' \times \mathbf{r})' = g(\mathbf{r}) \mathbf{r} \times \mathbf{r} = \mathbf{0}$ . Therefore,  $(\mathbf{r}' \times \mathbf{r}) = \mathbf{n}$ , a constant vector and so  $\mathbf{r} \cdot \mathbf{n} = \mathbf{r} \cdot (\mathbf{r}' \times \mathbf{r}) = 0$  showing that  $\mathbf{n}$  is a normal vector to a plane which contains  $\mathbf{r}(t)$  for all  $t$ . This proves the lemma.

The next example has as a special case one of Kepler's laws, Kepler's second law, the equal area law.



**Example 18.2.4** *An object moves in three dimensions in such a way that the only force acting on the object is a central force. Then the object moves in a plane and the radius vector from the origin to the object sweeps out area at a constant rate.*

The above lemma says the object moves in a plane. From the assumption that the force field is a central force field, it follows from (18.4) that

$$2r'(t)\theta'(t) + r(t)\theta''(t) = 0$$

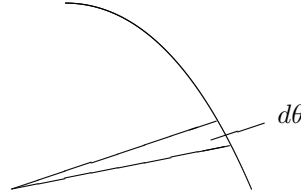
Multiply both sides of this equation by  $r$ . This yields

$$2rr'\theta' + r^2\theta'' = (r^2\theta')' = 0. \quad (18.5)$$

Consequently,

$$r^2\theta' = c \quad (18.6)$$

for some constant,  $C$ . Now consider the following picture.



In this picture,  $d\theta$  is the indicated angle and the two lines determining this angle are position vectors for the object at point  $t$  and point  $t + dt$ . The area of the circular sector,  $dA$ , is essentially  $r^2 d\theta$  and so  $dA = \frac{1}{2} r^2 d\theta$ . Therefore,

$$\frac{dA}{dt} = \frac{1}{2} r^2 \frac{d\theta}{dt} = \frac{c}{2}. \quad (18.7)$$

## 18.3 Planetary Motion

Kepler's laws of planetary motion state that planets move around the sun along an ellipse, the equal area law described above holds, and there is a formula for the time it takes for the planet to move around the sun. These laws, discovered by Kepler, were shown by Newton to be consequences of his law of gravitation which states that the force acting on a mass,  $m$  by a mass,  $M$  is given by

$$\mathbf{F} = -GMm \left( \frac{1}{r^3} \right) \mathbf{r} = -GMm \left( \frac{1}{r^2} \right) \mathbf{e}_r$$

where  $r$  is the distance between centers of mass and  $\mathbf{r}$  is the position vector from  $M$  to  $m$ . Here  $G$  is the gravitation constant. This is called an inverse square law. Gravity acts according to this law and so does electrostatic force. The constant,  $G$ , is very small when usual units are used and it has been computed using a very delicate experiment. It is now accepted to be

$$6.67 \times 10^{-11} \text{ Newton meter}^2/\text{kilogram}^2.$$

The experiment involved a light source shining on a mirror attached to a quartz fiber from which was suspended a long rod with two equal masses at the ends which were attracted

by two larger masses. The gravitation force between the suspended masses and the two large masses caused the fibre to twist ever so slightly and this twisting was measured by observing the deflection of the light reflected from the mirror on a scale placed some distance from the fibre. The constant was first measured successfully by Lord Cavendish in 1798 and the present accepted value was obtained in 1942. Experiments like these are major accomplishments.

In the following argument,  $M$  is the mass of the sun and  $m$  is the mass of the planet. (It could also be a comet or an asteroid.)

Consider the first of Kepler's laws, the one which states that planets move along ellipses. From Lemma 18.2.3, the motion is in a plane. Now from (18.4) and Newton's second law,

$$\left(r''(t) - r(t)\theta'(t)^2\right)\mathbf{e}_r + (2r'(t)\theta'(t) + r(t)\theta''(t))\mathbf{e}_\theta = -\frac{GMm}{m}\left(\frac{1}{r^2}\right)\mathbf{e}_r = -k\left(\frac{1}{r^2}\right)\mathbf{e}_r$$

Thus  $k = GM$  and

$$r''(t) - r(t)\theta'(t)^2 = -k\left(\frac{1}{r^2}\right), \quad 2r'(t)\theta'(t) + r(t)\theta''(t) = 0. \quad (18.8)$$

As in (18.5),  $(r^2\theta')' = 0$  and so there exists a constant,  $c$ , such that

$$r^2\theta' = c. \quad (18.9)$$

Therefore, also,

$$2rr'\theta' + r^2\theta'' = 0$$

and so

$$\theta'' = \frac{-2rr'\theta'}{r^2} = \frac{-2r'\theta'}{r} = \frac{-2}{r} \frac{dr}{dt} \frac{c}{r^2} \quad (18.10)$$

$$= \frac{-2c}{r^3} \frac{dr}{dt} \quad (18.11)$$

Now consider the first of the above equations. The question of interest is to know how  $r$  depends on  $\theta$ . By the chain rule, regarding  $r$  as a function of  $\theta$  and  $\theta$  as a function of  $t$ ,

$$\frac{dr}{d\theta} \frac{d\theta}{dt} = \frac{dr}{dt} \quad (18.12)$$

and by (18.9),

$$\frac{dr}{d\theta} = \frac{c}{r^2} \frac{dr}{dt}. \quad (18.13)$$

Also, by (18.10) and (18.9),

$$\begin{aligned} \frac{d^2\theta}{dt^2} &= \frac{-2c}{r^3} \frac{dr}{dt} = \frac{-2c}{r^3} \left(\frac{dr}{d\theta} \frac{d\theta}{dt}\right) \\ &= \frac{-2c}{r^3} \frac{dr}{d\theta} \left(\frac{c}{r^2}\right) = \frac{-2c^2}{r^5} \frac{dr}{d\theta} \end{aligned}$$

Differentiating (18.12) again with respect to  $t$ ,

$$\begin{aligned} \frac{d^2r}{dt^2} &= \frac{d^2r}{d\theta^2} \left(\frac{d\theta}{dt}\right)^2 + \frac{dr}{d\theta} \frac{d^2\theta}{dt^2} \\ &= \frac{d^2r}{d\theta^2} \left(\frac{c}{r^2}\right)^2 + \frac{dr}{d\theta} \left(\frac{-2c^2}{r^5} \frac{dr}{d\theta}\right) \\ &= \frac{d^2r}{d\theta^2} \left(\frac{c}{r^2}\right)^2 - \left(\frac{dr}{d\theta}\right)^2 \left(\frac{2c^2}{r^5}\right). \end{aligned}$$

It follows that the first equation of (18.8) yields  $r''(t) - r(t)\theta'(t)^2 = -k\left(\frac{1}{r^2}\right)$

$$\frac{d^2 r}{d\theta^2} \left(\frac{c}{r^2}\right)^2 - \left(\frac{dr}{d\theta}\right)^2 \left(\frac{2c^2}{r^5}\right) - r \left(\frac{c}{r^2}\right)^2 = -k \left(\frac{1}{r^2}\right)$$

which is a fairly messy looking differential equation. However, it can be simplified by multiplying both sides by  $\frac{r^4}{c^2}$  to get

$$\frac{d^2 r}{d\theta^2} - \left(\frac{dr}{d\theta}\right)^2 \left(\frac{2}{r}\right) - r = -\frac{k}{c^2} r^2 \quad (18.14)$$

Next consider the above equation in terms of  $\rho = r^{-1}$ . Thus, from the chain rule,

$$r = \rho^{-1}, \frac{dr}{d\theta} = (-1)\rho^{-2} \frac{d\rho}{d\theta},$$

$$\frac{d^2 r}{d\theta^2} = 2\rho^{-3} \left(\frac{d\rho}{d\theta}\right)^2 - \rho^{-2} \frac{d^2 \rho}{d\theta^2}.$$

substituting this in to (18.14),

$$\overbrace{2\rho^{-3} \left(\frac{d\rho}{d\theta}\right)^2 - \rho^{-2} \frac{d^2 \rho}{d\theta^2}}^{\frac{d^2 r}{d\theta^2}} - \left(\overbrace{(-1)\rho^{-2} \frac{d\rho}{d\theta}}^{\frac{dr}{d\theta}}\right)^2 (2\rho) - \rho^{-1} = -\frac{k}{\rho^2 c^2}.$$

Now note that the first and third terms add to zero and so

$$-\rho^{-2} \frac{d^2 \rho}{d\theta^2} - \rho^{-1} = -\frac{k}{\rho^2 c^2}$$

Multiplying both sides by  $-\rho^{-2}$  yields the equation,

$$\frac{d^2 \rho}{d\theta^2} + \rho = \frac{k}{c^2},$$

a much more manageable equation. Multiply both sides by  $\frac{d\rho}{d\theta}$ .

$$\frac{d^2 \rho}{d\theta^2} \frac{d\rho}{d\theta} + \rho \frac{d\rho}{d\theta} = \frac{k}{c^2} \frac{d\rho}{d\theta}$$

Then from the product rule,

$$\frac{1}{2} \frac{d}{d\theta} \left( \left(\frac{d\rho}{d\theta}\right)^2 + \rho^2 \right) = \frac{k}{c^2} \frac{d\rho}{d\theta}.$$

Therefore, there exists a constant,  $c_1$  such that

$$\frac{1}{2} \left( \left(\frac{d\rho}{d\theta}\right)^2 + \rho^2 \right) - \frac{k}{c^2} \rho = c_1$$

and so

$$\begin{aligned} \left(\frac{d\rho}{d\theta}\right)^2 &= 2c_1 + \frac{k}{c^2}\rho - \rho^2 \\ &= \left(\frac{k^2}{4c^4} + 2c_1\right) - \left(\rho - \frac{k}{2c^2}\right)^2 \\ &\equiv \delta^2 - \left(\rho - \frac{k}{2c^2}\right)^2 \end{aligned}$$

Now letting  $\rho_1 = \rho - \frac{k}{2c^2}$ ,

$$\frac{1}{\delta^2} \left(\frac{d\rho_1}{d\theta}\right)^2 + \left(\frac{\rho_1}{\delta}\right)^2 = 1$$

which shows that  $\left(\frac{1}{\delta} \frac{d\rho_1}{d\theta}, \frac{\rho_1}{\delta}\right)$  is a point on the unit circle. Therefore, there exists an angle,  $\alpha(\theta)$  such that

$$\frac{d\rho_1}{d\theta} = \delta \cos(\alpha(\theta)), \rho_1 = \delta \sin(\alpha(\theta)).$$

Differentiating the second equation with respect to  $\theta$ ,

$$\frac{d\rho_1}{d\theta} = \alpha'(\theta) \delta \cos(\alpha(\theta))$$

and so  $\alpha'(\theta) = 1$ . Therefore,  $\alpha(\theta) = \theta + \phi$  for some constant,  $\phi$ . Redefining,  $\theta$  if necessary, (Let  $\tilde{\theta} = \theta + \phi$ ) it can be assumed that  $\phi = 0$  so

$$\rho - \frac{k}{2c^2} = \rho_1 = \delta \sin \theta.$$

Thus

$$\rho = \frac{k}{c^2} + \delta \sin \theta$$

and so

$$\begin{aligned} r &= \frac{1}{\frac{k}{c^2} + \delta \sin \theta} = \frac{c^2/k}{1 + (c^2/k) \delta \sin \theta} \\ &= \frac{p\varepsilon}{1 + \varepsilon \sin \theta} \end{aligned} \tag{18.15}$$

where

$$\varepsilon = (c^2/k) \delta \text{ and } p = c^2/k\varepsilon. \tag{18.16}$$

Here all these constants are nonnegative.

Thus

$$r + \varepsilon r \sin \theta = \varepsilon p$$

and so  $r = (\varepsilon p - \varepsilon y)$ . Then squaring both sides,

$$x^2 + y^2 = (\varepsilon p - \varepsilon y)^2 = \varepsilon^2 p^2 - 2p\varepsilon^2 y + \varepsilon^2 y^2$$

And so

$$x^2 + (1 - \varepsilon^2) y^2 = \varepsilon^2 p^2 - 2p\varepsilon^2 y. \tag{18.17}$$

In case  $\varepsilon = 1$ , this reduces to the equation of a parabola. If  $\varepsilon < 1$ , this reduces to the equation of an ellipse and if  $\varepsilon > 1$ , this is called a hyperbola. This proves that objects which

are acted on only by a force of the form given in the above example move along hyperbolas, ellipses or circles. The case where  $\varepsilon = 0$  corresponds to a circle. The constant,  $\varepsilon$  is called the eccentricity. This is called Kepler's first law in the case of a planet.

Kepler's third law involves the time it takes for the planet to orbit the sun. From (18.17) you can complete the square and obtain

$$x^2 + (1 - \varepsilon^2) \left( y + \frac{p\varepsilon^2}{1 - \varepsilon^2} \right)^2 = \varepsilon^2 p^2 + \frac{p^2 \varepsilon^4}{(1 - \varepsilon^2)} = \frac{\varepsilon^2 p^2}{(1 - \varepsilon^2)},$$

and this yields

$$x^2 / \left( \frac{\varepsilon^2 p^2}{1 - \varepsilon^2} \right) + \left( y + \frac{p\varepsilon^2}{1 - \varepsilon^2} \right)^2 / \left( \frac{\varepsilon^2 p^2}{(1 - \varepsilon^2)^2} \right) = 1. \quad (18.18)$$

Now note this is the equation of an ellipse and that the diameter of this ellipse is

$$\frac{2\varepsilon p}{(1 - \varepsilon^2)} \equiv 2a.$$

This follows because

$$\frac{\varepsilon^2 p^2}{(1 - \varepsilon^2)^2} \geq \frac{\varepsilon^2 p^2}{1 - \varepsilon^2}.$$

Now let  $T$  denote the time it takes for the planet to make one revolution about the sun. Recall Example 10.9.2 on Page 248 which gives the area of an ellipse. Using this formula, and (18.7) the following equation must hold.

$$\overbrace{\pi \frac{\varepsilon p}{\sqrt{1 - \varepsilon^2}} \frac{\varepsilon p}{(1 - \varepsilon^2)}}^{\text{area of ellipse}} = T \frac{c}{2}$$

Therefore,

$$T = \frac{2}{c} \frac{\pi \varepsilon^2 p^2}{(1 - \varepsilon^2)^{3/2}}$$

and so

$$T^2 = \frac{4\pi^2 \varepsilon^4 p^4}{c^2 (1 - \varepsilon^2)^3}$$

Now using (18.16),

$$\begin{aligned} T^2 &= \frac{4\pi^2 \varepsilon^4 p^4}{k \varepsilon p (1 - \varepsilon^2)^3} = \frac{4\pi^2 (\varepsilon p)^3}{k (1 - \varepsilon^2)^3} \\ &= \frac{4\pi^2 a^3}{k} = \frac{4\pi^2 a^3}{GM}. \end{aligned}$$

Written more memorably, this has shown

$$T^2 = \frac{4\pi^2}{GM} \left( \frac{\text{diameter of ellipse}}{2} \right)^3. \quad (18.19)$$

This relationship is known as Kepler's third law. Kepler's second law, the equal area formula, holds for any central force, not just one which satisfies an inverse square law.

## 18.4 Exercises

1. In general it is a stupid idea to try to use algebra to invert and solve for a set of curvilinear coordinates such as polar or cylindrical coordinates in term of Cartesian coordinates. Not only is it often very difficult or even impossible to do it, but also it takes you in entirely the wrong direction because the whole point of introducing the new coordinates is to write everything in terms of these new coordinates and not in terms of Cartesian coordinates. However, sometimes this inversion can be done. Describe how to solve for  $r$  and  $\theta$  in terms of  $x$  and  $y$  in polar coordinates.
2. A point has Cartesian coordinates,  $(1, 2, 3)$ . Find its spherical and cylindrical coordinates.
3. Describe the following surface in rectangular coordinates.  $\phi = \pi/4$  where  $\phi$  is the polar angle in spherical coordinates.
4. Describe the following surface in rectangular coordinates.  $\theta = \pi/4$  where  $\theta$  is the angle measured from the positive  $x$  axis spherical coordinates.
5. Describe the following surface in rectangular coordinates.  $\theta = \pi/4$  where  $\theta$  is the angle measured from the positive  $x$  axis cylindrical coordinates.
6. Describe the following surface in rectangular coordinates.  $r = 5$  where  $r$  is one of the cylindrical coordinates.
7. Describe the following surface in rectangular coordinates.  $\rho = 4$  where  $\rho$  is the distance to the origin.
8. Give the cone,  $z = \sqrt{x^2 + y^2}$  in cylindrical coordinates and in spherical coordinates.
9. Write the following in spherical coordinates.
  - (a)  $z = x^2 + y^2$ .
  - (b)  $x^2 - y^2 = 1$
  - (c)  $z^2 + x^2 + y^2 = 6$
  - (d)  $z = \sqrt{x^2 + y^2}$
  - (e)  $y = x$
  - (f)  $z = x$
10. Write the following in cylindrical coordinates.
  - (a)  $z = x^2 + y^2$ .
  - (b)  $x^2 - y^2 = 1$
  - (c)  $z^2 + x^2 + y^2 = 6$
  - (d)  $z = \sqrt{x^2 + y^2}$
  - (e)  $y = x$
  - (f)  $z = x$
11. Suppose you know how the spherical coordinates of a moving point change as a function of  $t$ . Can you figure out the velocity of the point? Specifically, suppose  $\phi(t) = t$ ,  $\theta(t) = 1 + t$ , and  $\rho(t) = t$ . Find the speed and the velocity of the object in terms of Cartesian coordinates. **Hint:** You would need to find  $x'(t)$ ,  $y'(t)$ , and  $z'(t)$ . Then in terms of Cartesian coordinates, the velocity would be  $x'(t)\mathbf{i} + y'(t)\mathbf{j} + z'(t)\mathbf{k}$ .

12. Explain why low pressure areas rotate counter clockwise in the Northern hemisphere and clockwise in the Southern hemisphere. **Hint:** Note that from the point of view of an observer fixed in space, the low pressure area already has a counter clockwise rotation because of the rotation of the earth and its spherical shape. Now consider (18.6). In the low pressure area stuff will move toward the center so  $r$  gets smaller. How are things different in the Southern hemisphere?
13. What are some physical assumptions which are made in the above derivation of Keplers laws from Newton's laws of motion?
14. The orbit of the earth is pretty nearly circular and the distance from the sun to the earth is about  $149 \times 10^6$  kilometers. Using (18.19) and the above value of the universal gravitation constant, determine the mass of the sun. The earth goes around it in 365 days. (Actually it is 365.256 days.)
15. It is desired to place a satellite above the equator of the earth which will rotate about the center of mass of the earth every 24 hours. Is it necessary that the orbit be circular? What if you want the satellite to stay above the same point on the earth at all times? If the orbit is to be circular and the satellite is to stay above the same point, at what distance from the center of mass of the earth should the satellite be? You may use that the mass of the earth is  $5.98 \times 10^{24}$  kilograms. Such a satellite is called geosynchronous.





## Part IV

# Functions Of More Than One Variable



# Linear Algebra

## 19.1 Matrices

A matrix is a rectangular array of numbers. Several of them are referred to as matrices. For example, here is a matrix.

$$\begin{pmatrix} 1 & 2 & 3 & 4 \\ 5 & 2 & 8 & 7 \\ 6 & -9 & 1 & 2 \end{pmatrix}$$

This matrix is a  $3 \times 4$  matrix because there are three rows and four columns. The first row is  $(1\ 2\ 3\ 4)$ , the second row is  $(5\ 2\ 8\ 7)$  and so forth. The first column is  $\begin{pmatrix} 1 \\ 5 \\ 6 \end{pmatrix}$ . The

convention in dealing with matrices is to always list the rows first and then the columns. Also, you can remember the columns are like columns in a Greek temple. They stand up right while the rows just lay there like rows made by a tractor in a plowed field. Elements of the matrix are identified according to position in the matrix. For example, 8 is in position 2,3 because it is in the second row and the third column. The symbol,  $(a_{ij})$  refers to a matrix in which the  $i$  denotes the row and the  $j$  denotes the column. Using this notation on the above matrix,  $a_{23} = 8$ ,  $a_{32} = -9$ ,  $a_{12} = 2$ , etc.

There are various operations which are done on matrices. They can sometimes be added, multiplied by a scalar and sometimes multiplied. To illustrate scalar multiplication, consider the following example.

$$3 \begin{pmatrix} 1 & 2 & 3 & 4 \\ 5 & 2 & 8 & 7 \\ 6 & -9 & 1 & 2 \end{pmatrix} = \begin{pmatrix} 3 & 6 & 9 & 12 \\ 15 & 6 & 24 & 21 \\ 18 & -27 & 3 & 6 \end{pmatrix}.$$

The new matrix is obtained by multiplying every entry of the original matrix by the given scalar. If  $A$  is an  $m \times n$  matrix,  $-A$  is defined to equal  $(-1)A$ .

Two matrices which are the same size can be added. When this is done, the result is the matrix which is obtained by adding corresponding entries. Thus

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 2 \end{pmatrix} + \begin{pmatrix} -1 & 4 \\ 2 & 8 \\ 6 & -4 \end{pmatrix} = \begin{pmatrix} 0 & 6 \\ 5 & 12 \\ 11 & -2 \end{pmatrix}.$$

Two matrices are equal exactly when they are the same size and the corresponding entries are identical.

Stated in a general form the rules for scalar multiplication and addition are these.

**Definition 19.1.1** Let  $A = (a_{ij})$  and  $B = (b_{ij})$  be two  $m \times n$  matrices. Then  $A + B = C$  where

$$C = (c_{ij})$$

for  $c_{ij} = a_{ij} + b_{ij}$ . Also if  $x$  is a scalar,

$$xA = (c_{ij})$$

where  $c_{ij} = xa_{ij}$ . The number  $A_{ij}$  will typically refer to the  $ij^{\text{th}}$  entry of the matrix,  $A$ . The zero matrix, denoted by  $0$  will be the matrix consisting of all zeros.

Note there are  $2 \times 3$  zero matrices,  $3 \times 4$  zero matrices, etc. In fact for every size there is a zero matrix.

With this definition, the following properties are all obvious but you should verify all of these properties valid for  $A$ ,  $B$ , and  $C$ ,  $m \times n$  matrices and  $0$  an  $m \times n$  zero matrix,

$$A + B = B + A, \quad (19.1)$$

the commutative law of addition,

$$(A + B) + C = A + (B + C), \quad (19.2)$$

the associative law for addition,

$$A + 0 = A, \quad (19.3)$$

the existence of an additive identity,

$$A + (-A) = 0, \quad (19.4)$$

the existence of an additive inverse. Also, for  $\alpha, \beta$  scalars, the following also hold.

$$\alpha(A + B) = \alpha A + \alpha B, \quad (19.5)$$

$$(\alpha + \beta)A = \alpha A + \beta A, \quad (19.6)$$

$$\alpha(\beta A) = \alpha\beta(A), \quad (19.7)$$

$$1A = A. \quad (19.8)$$

The above properties, (19.1) - (19.8) are known as the vector space axioms and the fact that the  $m \times n$  matrices satisfy these axioms is what is meant by saying this set of matrices forms a vector space.

**Definition 19.1.2** Let  $\mathbf{x} \in \mathbb{R}^n$ . Recall this means  $\mathbf{x}$  is a list of  $n$  real numbers. From now on, this list of numbers will be arranged vertically. Thus

$$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}.$$

These  $n \times 1$  matrices are also called vectors, sometimes column vectors, while a  $1 \times n$  matrix of the form  $(x_1 \cdots x_n)$  is referred to as a row vector.

All the above is fine, but the real reason for considering matrices is that they can be multiplied. This is where things quit being banal.

First consider the problem of multiplying an  $m \times n$  matrix by an  $n \times 1$  column vector. Consider the following example

$$\begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} 7 \\ 8 \\ 9 \end{pmatrix} = ?$$

The way I like to remember this is as follows. Slide the vector, placing it on top the two rows as shown

$$\begin{pmatrix} 7 & 8 & 9 \\ 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix},$$

multiply the numbers on the top by the numbers on the bottom and add them up to get a single number for each row of the matrix. These numbers are listed in the same order giving, in this case, a  $2 \times 1$  matrix. Thus

$$\begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} 7 \\ 8 \\ 9 \end{pmatrix} = \begin{pmatrix} 7 \times 1 + 8 \times 2 + 9 \times 3 \\ 7 \times 4 + 8 \times 5 + 9 \times 6 \end{pmatrix} = \begin{pmatrix} 50 \\ 122 \end{pmatrix}.$$

In more general terms,

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 \end{pmatrix}.$$

Motivated by this example, here is the definition of how to multiply an  $m \times n$  matrix by an  $n \times 1$  matrix. (vector)

**Definition 19.1.3** Let  $A = A_{ij}$  be an  $m \times n$  matrix and let  $\mathbf{v}$  be an  $n \times 1$  matrix,

$$\mathbf{v} = \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix}$$

Then  $A\mathbf{v}$  is an  $m \times 1$  matrix and the  $i^{\text{th}}$  component of this matrix is

$$(A\mathbf{v})_i = \sum_{j=1}^n A_{ij}v_j.$$

Thus

$$A\mathbf{v} = \begin{pmatrix} \sum_{j=1}^n A_{1j}v_j \\ \vdots \\ \sum_{j=1}^n A_{mj}v_j \end{pmatrix}.$$

Using the repeated index summation convention on Page 331,

$$(A\mathbf{v})_i = A_{ij}v_j.$$

Note also that multiplication by an  $m \times n$  matrix takes a vector in  $\mathbb{R}^n$ , an  $n \times 1$  matrix, and produces a vector in  $\mathbb{R}^m$ , an  $m \times 1$  matrix.

With this done, the next task is to multiply an  $m \times n$  matrix times an  $n \times p$  matrix. Before doing so, the following may be helpful.

$$\widehat{(m \times n)(n \times p)}^{\text{these must match}} = m \times p$$

**If the two middle numbers don't match, you can't multiply the matrices.**

Let  $A$  be an  $m \times n$  matrix and let  $B$  be an  $n \times p$  matrix. Then  $B$  is of the form

$$B = (\mathbf{b}_1, \dots, \mathbf{b}_p)$$

where  $\mathbf{b}_k$  is an  $n \times 1$  matrix. Then an  $m \times p$  matrix,  $AB$  is defined as follows:

$$AB \equiv (A\mathbf{b}_1, \dots, A\mathbf{b}_p) \quad (19.9)$$

where  $A\mathbf{b}_k$  is an  $m \times 1$  matrix. Hence  $AB$  as just defined is an  $m \times p$  matrix.

What is the  $ij^{th}$  entry of  $AB$ ? It would be the  $i^{th}$  entry of the  $j^{th}$  column of  $AB$ . Thus it would be the  $i^{th}$  entry of  $A\mathbf{b}_j$ . Now

$$\mathbf{b}_j = \begin{pmatrix} B_{1j} \\ \vdots \\ B_{nj} \end{pmatrix}$$

and from the above definition, the  $i^{th}$  entry is

$$\sum_{k=1}^n A_{ik} B_{kj}. \quad (19.10)$$

This motivates the definition for matrix multiplication which identifies the  $ij^{th}$  entries of the product.

**Definition 19.1.4** Let  $A = (A_{ij})$  be an  $m \times n$  matrix and let  $B = (B_{ij})$  be an  $n \times p$  matrix. Then  $AB$  is an  $m \times p$  matrix and

$$(AB)_{ij} = \sum_{k=1}^n A_{ik} B_{kj}. \quad (19.11)$$

In terms of the repeated summation convention, on Page 331,

$$(AB)_{ij} = A_{ik} B_{kj}. \quad (19.12)$$

Two matrices,  $A$  and  $B$  are said to be conformable in a particular order if they can be multiplied in that order. Thus if  $A$  is an  $r \times s$  matrix and  $B$  is a  $s \times p$  then  $A$  and  $B$  are conformable in the order,  $AB$ .

**Example 19.1.5** Multiply if possible  $\begin{pmatrix} 1 & 2 \\ 3 & 1 \\ 2 & 6 \end{pmatrix} \begin{pmatrix} 2 & 3 & 1 \\ 7 & 6 & 2 \end{pmatrix}$ .

First check to see if this is possible. It is of the form  $(3 \times 2)(2 \times 3)$  and since the inside numbers match, it must be possible to do this. The answer is of the form

$$\left( \begin{pmatrix} 1 & 2 \\ 3 & 1 \\ 2 & 6 \end{pmatrix} \begin{pmatrix} 2 \\ 7 \end{pmatrix}, \begin{pmatrix} 1 & 2 \\ 3 & 1 \\ 2 & 6 \end{pmatrix} \begin{pmatrix} 3 \\ 6 \end{pmatrix}, \begin{pmatrix} 1 & 2 \\ 3 & 1 \\ 2 & 6 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} \right)$$

where the commas separate the columns in the resulting product. Thus the above product equals

$$\begin{pmatrix} 16 & 15 & 5 \\ 13 & 15 & 5 \\ 46 & 42 & 14 \end{pmatrix},$$

a  $3 \times 3$  matrix as desired.

**Example 19.1.6** *Multiply if possible*  $\begin{pmatrix} 1 & 2 \\ 3 & 1 \\ 2 & 6 \end{pmatrix} \begin{pmatrix} 2 & 3 & 1 \\ 7 & 6 & 2 \\ 0 & 0 & 0 \end{pmatrix}$ .

This is not possible because it is of the form  $(3 \times 2)(3 \times 3)$  and the middle numbers don't match.

**Example 19.1.7** *Multiply if possible*  $\begin{pmatrix} 2 & 3 & 1 \\ 7 & 6 & 2 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 3 & 1 \\ 2 & 6 \end{pmatrix}$ .

This is possible because in this case it is of the form  $(3 \times 3)(3 \times 2)$  and the middle numbers do match. When the multiplication is done it equals

$$\begin{pmatrix} 13 & 13 \\ 29 & 32 \\ 0 & 0 \end{pmatrix}.$$

Note the above two examples show matrix multiplication is not commutative. In fact, it can happen that  $AB$  makes sense while  $BA$  does not make sense which was the case in these examples. It seems natural to ask if the product of any two  $n \times n$  matrices commutes.

**Example 19.1.8** *Compare*  $\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$  and  $\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$ .

The first product is

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 2 & 1 \\ 4 & 3 \end{pmatrix},$$

the second product is

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} = \begin{pmatrix} 3 & 4 \\ 1 & 2 \end{pmatrix},$$

and you see these are not equal. Therefore, you cannot conclude that  $AB = BA$  for matrix multiplication. However, there are some properties which do hold.

**Proposition 19.1.9** *If all multiplications and additions make sense, the following hold for matrices and  $a, b$  scalars.*

$$A(aB + bC) = a(AB) + b(AC) \tag{19.13}$$

$$(B + C)A = BA + CA \tag{19.14}$$

$$A(BC) = (AB)C \tag{19.15}$$

**Proof:** Using the repeated index summation convention and the above definition of matrix multiplication,

$$\begin{aligned}
 (A(aB + bC))_{ij} &= A_{ik}(aB + bC)_{kj} \\
 &= A_{ik}(aB_{kj} + bC_{kj}) \\
 &= aA_{ik}B_{kj} + bA_{ik}C_{kj} \\
 &= a(AB)_{ij} + b(AC)_{ij} \\
 &= (a(AB) + b(AC))_{ij}
 \end{aligned}$$

showing that  $A(B + C) = AB + AC$  as claimed. Formula (19.14) is entirely similar.

Consider (19.15), the associative law of multiplication.

$$\begin{aligned}
 (A(BC))_{ij} &= A_{ik}(BC)_{kj} \\
 &= A_{ik}B_{kl}C_{lj} \\
 &= (AB)_{il}C_{lj} \\
 &= ((AB)C)_{ij}.
 \end{aligned}$$

This proves (19.15).

Another important operation on matrices is that of taking the transpose. The following example shows what is meant by this operation, denoted by placing a  $T$  as an exponent on the matrix.

$$\begin{pmatrix} 1 & 2 \\ 3 & 1 \\ 2 & 6 \end{pmatrix}^T = \begin{pmatrix} 1 & 3 & 2 \\ 2 & 1 & 6 \end{pmatrix}$$

What happened? The first column became the first row and the second column became the second row. Thus the  $3 \times 2$  matrix became a  $2 \times 3$  matrix. The number 3 was in the second row and the first column and it ended up in the first row and second column. This motivates the following definition of the transpose of a matrix.

**Definition 19.1.10** Let  $A$  be an  $m \times n$  matrix. Then  $A^T$  denotes the  $n \times m$  matrix which is defined as follows.

$$(A^T)_{ij} = A_{ji}$$

The transpose of a matrix has the following important property.

**Lemma 19.1.11** Let  $A$  be an  $m \times n$  matrix and let  $B$  be a  $n \times p$  matrix. Then

$$(AB)^T = B^T A^T \quad (19.16)$$

and if  $\alpha$  and  $\beta$  are scalars,

$$(\alpha A + \beta B)^T = \alpha A^T + \beta B^T \quad (19.17)$$

**Proof:** From the definition,

$$\begin{aligned}
 ((AB)^T)_{ij} &= (AB)_{ji} \\
 &= A_{jk}B_{ki} \\
 &= (B^T)_{ik}(A^T)_{kj} \\
 &= (B^T A^T)_{ij}
 \end{aligned}$$

(19.17) is left as an exercise and this proves the lemma.



There is a special matrix called  $I$  and defined by

$$I_{ij} = \delta_{ij}$$

where  $\delta_{ij}$  is the Kroneker symbol defined on Page 330. It is called the identity matrix because it is a multiplicative identity in the following sense.

**Lemma 19.1.12** *Suppose  $A$  is an  $m \times n$  matrix and  $I_n$  is the  $n \times n$  identity matrix. Then  $AI_n = A$ . If  $I_m$  is the  $m \times m$  identity matrix, it also follows that  $I_mA = A$ .*

**Proof:**

$$\begin{aligned}(AI_n)_{ij} &= A_{ik}\delta_{kj} \\ &= A_{ij}\end{aligned}$$

and so  $AI_n = A$ . The other case is left as an exercise for you.

**Definition 19.1.13** *An  $n \times n$  matrix,  $A$  has an inverse,  $A^{-1}$  if and only if  $AA^{-1} = A^{-1}A = I$  where  $I = (\delta_{ij})$  for*

$$\delta_{ij} \equiv \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

### 19.1.1 Finding The Inverse Of A Matrix

A little later a formula is given for the inverse of a matrix. However, it is not a good way to find the inverse for a matrix. It is also important to note that not all matrices have inverses.

**Example 19.1.14** *Let  $A = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$ . Does  $A$  have an inverse?*

One might think  $A$  would have an inverse because it does not equal zero. However,

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} -1 & \\ & 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

and if  $A^{-1}$  existed, this could not happen because you could write

$$\begin{aligned}\begin{pmatrix} 0 \\ 0 \end{pmatrix} &= A^{-1} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix} \right) = A^{-1} \left( A \begin{pmatrix} -1 \\ 1 \end{pmatrix} \right) = \\ &= (A^{-1}A) \begin{pmatrix} -1 \\ 1 \end{pmatrix} = I \begin{pmatrix} -1 \\ 1 \end{pmatrix} = \begin{pmatrix} -1 \\ 1 \end{pmatrix},\end{aligned}$$

a contradiction. Thus the answer is that  $A$  does not have an inverse.

**Example 19.1.15** *Let  $A = \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix}$ . Show  $\begin{pmatrix} 2 & -1 \\ -1 & 1 \end{pmatrix}$  is the inverse of  $A$ .*

To check this, multiply

$$\begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} 2 & -1 \\ -1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

and

$$\begin{pmatrix} 2 & -1 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

showing that this matrix is indeed the inverse of  $A$ .

In the last example, how would you find  $A^{-1}$ ? You wish to find a matrix,  $\begin{pmatrix} x & z \\ y & w \end{pmatrix}$  such that

$$\begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} x & z \\ y & w \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}.$$

This requires the solution of the systems of equations,

$$x + y = 1, x + 2y = 0$$

and

$$z + w = 0, z + 2w = 1.$$

Writing the augmented matrix for these two systems gives

$$\begin{pmatrix} 1 & 1 & 1 \\ 1 & 2 & 0 \end{pmatrix} \quad (19.18)$$

for the first system and

$$\begin{pmatrix} 1 & 1 & 0 \\ 1 & 2 & 1 \end{pmatrix} \quad (19.19)$$

for the second. Lets solve the first system. Take  $(-1)$  times the first row and add to the second to get

$$\begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \end{pmatrix}$$

Now take  $(-1)$  times the second row and add to the first to get

$$\begin{pmatrix} 1 & 0 & 2 \\ 0 & 1 & -1 \end{pmatrix}.$$

Putting in the variables, this says  $x = 2$  and  $y = -1$ .

Now solve the second system, (19.19) to find  $z$  and  $w$ . Take  $(-1)$  times the first row and add to the second to get

$$\begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \end{pmatrix}.$$

Now take  $(-1)$  times the second row and add to the first to get

$$\begin{pmatrix} 1 & 0 & -1 \\ 0 & 1 & 1 \end{pmatrix}.$$

Putting in the variables, this says  $z = -1$  and  $w = 1$ . Therefore, the inverse is

$$\begin{pmatrix} 2 & -1 \\ -1 & 1 \end{pmatrix}.$$

Didn't the above seem rather repetitive? Note that exactly the same row operations were used in both systems. In each case, the end result was something of the form  $(I|\mathbf{v})$  where  $I$  is the identity and  $\mathbf{v}$  gave a column of the inverse. In the above,  $\begin{pmatrix} x \\ y \end{pmatrix}$ , the first column of the inverse was obtained first and then the second column  $\begin{pmatrix} z \\ w \end{pmatrix}$ .

This is the reason for the following simple procedure for finding the inverse of a matrix.

**Procedure 19.1.16** Suppose  $A$  is an  $n \times n$  matrix. To find  $A^{-1}$  if it exists, form the augmented  $n \times 2n$  matrix,

$$(A|I)$$

and then do row operations until you obtain an  $n \times 2n$  matrix of the form

$$(I|B) \tag{19.20}$$

if possible. When this has been done,  $B = A^{-1}$ . The matrix,  $A$  has no inverse exactly when it is impossible to do row operations and end up with one like (19.20).

**Example 19.1.17** Let  $A = \begin{pmatrix} 1 & 0 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & -1 \end{pmatrix}$ . Find  $A^{-1}$ .

Form the augmented matrix,

$$\begin{pmatrix} 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & -1 & 1 & 0 & 1 & 0 \\ 1 & 1 & -1 & 0 & 0 & 1 \end{pmatrix}.$$

Now do row operations until the  $n \times n$  matrix on the left becomes the identity matrix. This yields after a some computations,

$$\begin{pmatrix} 1 & 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 1 & 0 & 1 & -1 & 0 \\ 0 & 0 & 1 & 1 & -\frac{1}{2} & -\frac{1}{2} \end{pmatrix}$$

and so the inverse of  $A$  is the matrix on the right,

$$\begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & -1 & 0 \\ 1 & -\frac{1}{2} & -\frac{1}{2} \end{pmatrix}.$$

Checking the answer is easy. Just multiply the matrices and see if it works.

$$\begin{pmatrix} 1 & 0 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & -1 \end{pmatrix} \begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & -1 & 0 \\ 1 & -\frac{1}{2} & -\frac{1}{2} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

## 19.2 Exercises

1. In (19.1) - (19.8) describe  $-A$  and  $0$ .
2. Using only the properties (19.1) - (19.8) show  $-A$  is unique.
3. Using only the properties (19.1) - (19.8) show  $0$  is unique.
4. Using only the properties (19.1) - (19.8) show  $0A = 0$ . Here the  $0$  on the left is the scalar  $0$  and the  $0$  on the right is the zero for  $m \times n$  matrices.
5. Using only the properties (19.1) - (19.8) and previous problems show  $(-1)A = -A$ .
6. Prove (19.17).
7. Prove that  $I_m A = A$  where  $A$  is an  $m \times n$  matrix.

8. Let  $A$  be a real  $m \times n$  matrix and let  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{y} \in \mathbb{R}^m$ . Show  $(A\mathbf{x}, \mathbf{y})_{\mathbb{R}^m} = (\mathbf{x}, A^T\mathbf{y})_{\mathbb{R}^n}$  where  $(\cdot, \cdot)_{\mathbb{R}^k}$  denotes the dot product in  $\mathbb{R}^k$ .
9. Use the result of Problem 8 to verify directly that  $(AB)^T = B^T A^T$  without making any reference to subscripts.
10. Let  $\mathbf{x} = (-1, -1, 1)$  and  $\mathbf{y} = (0, 1, 2)$ . Find  $\mathbf{x}^T \mathbf{y}$  and  $\mathbf{x} \mathbf{y}^T$  if possible.
11. Give an example of matrices,  $A, B, C$  such that  $B \neq C$ ,  $A \neq 0$ , and yet  $AB = AC$ .
12. Let  $A = \begin{pmatrix} 1 & 1 \\ -2 & -1 \\ 1 & 2 \end{pmatrix}$ ,  $B = \begin{pmatrix} 1 & -1 & -2 \\ 2 & 1 & -2 \end{pmatrix}$ , and  $C = \begin{pmatrix} 1 & 1 & -3 \\ -1 & 2 & 0 \\ -3 & -1 & 0 \end{pmatrix}$ . Find if possible.
  - (a)  $AB$
  - (b)  $BA$
  - (c)  $AC$
  - (d)  $CA$
  - (e)  $CB$
  - (f)  $BC$
13. Show that if  $A^{-1}$  exists for an  $n \times n$  matrix, then it is unique. That is, if  $BA = I$  and  $AB = I$ , then  $B = A^{-1}$ .
14. Show  $(AB)^{-1} = B^{-1}A^{-1}$ .
15. Give an example of a matrix,  $A$  such that  $A^2 = I$  and yet  $A \neq I$  and  $A \neq -I$ .
16. Give an example of matrices,  $A, B$  such that neither  $A$  nor  $B$  equals zero and yet  $AB = 0$ .
17. Write  $\begin{pmatrix} x_1 - x_2 + 2x_3 \\ 2x_3 + x_1 \\ 3x_3 \\ 3x_4 + 3x_2 + x_1 \end{pmatrix}$  in the form  $A \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix}$  where  $A$  is an appropriate matrix.
18. Give another example other than the one given in this section of two square matrices,  $A$  and  $B$  such that  $AB \neq BA$ .
19. Suppose  $A$  and  $B$  are square matrices of the same size. Which of the following are correct?
  - (a)  $(A - B)^2 = A^2 - 2AB + B^2$
  - (b)  $(AB)^2 = A^2B^2$
  - (c)  $(A + B)^2 = A^2 + 2AB + B^2$
  - (d)  $(A + B)^2 = A^2 + AB + BA + B^2$
  - (e)  $A^2B^2 = A(AB)B$
  - (f)  $(A + B)^3 = A^3 + 3A^2B + 3AB^2 + B^3$
  - (g)  $(A + B)(A - B) = A^2 - B^2$
  - (h) None of the above. They are all wrong.

(i) All of the above. They are all right.

20. Let  $A = \begin{pmatrix} -1 & -1 \\ 3 & 3 \end{pmatrix}$ . Find all  $2 \times 2$  matrices,  $B$  such that  $AB = 0$ .

21. Prove that if  $A^{-1}$  exists and  $A\mathbf{x} = \mathbf{0}$  then  $\mathbf{x} = \mathbf{0}$ .

22. Let

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 1 & 4 \\ 1 & 0 & 2 \end{pmatrix}.$$

Find  $A^{-1}$  if possible. If  $A^{-1}$  does not exist, determine why.

### 19.3 Linear Transformations

By (19.13), if  $A$  is an  $m \times n$  matrix, then for  $\mathbf{v}, \mathbf{u}$  vectors in  $\mathbb{R}^n$  and  $a, b$  scalars,

$$A \left( \overbrace{a\mathbf{u} + b\mathbf{v}}^{\in \mathbb{R}^n} \right) = aA\mathbf{u} + bA\mathbf{v} \in \mathbb{R}^m \quad (19.21)$$

**Definition 19.3.1** A function,  $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is called a linear transformation if for all  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$  and  $a, b$  scalars, (19.21) holds.

From (19.21), matrix multiplication defines a linear transformation as just defined. It turns out this is the only type of linear transformation available. Thus if  $A$  is a linear transformation from  $\mathbb{R}^n$  to  $\mathbb{R}^m$ , there is always a matrix which produces  $A$ . Before showing this, here is a simple definition.

**Definition 19.3.2** A vector,  $\mathbf{e}_i \in \mathbb{R}^n$  is defined as follows:

$$\mathbf{e}_i \equiv \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix},$$

where the 1 is in the  $i^{\text{th}}$  position and there are zeros everywhere else. Thus

$$\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)^T.$$

Of course the  $\mathbf{e}_i$  for a particular value of  $i$  in  $\mathbb{R}^n$  would be different than the  $\mathbf{e}_i$  for that same value of  $i$  in  $\mathbb{R}^m$  for  $m \neq n$ . One of them is longer than the other. However, which one is meant will be determined by the context in which they occur.

These vectors have a significant property.

**Lemma 19.3.3** Let  $\mathbf{v} \in \mathbb{R}^n$ . Thus  $\mathbf{v}$  is a list of real numbers arranged vertically,  $v_1, \dots, v_n$ . Then

$$\mathbf{e}_i^T \mathbf{v} = v_i. \quad (19.22)$$

Also, if  $A$  is an  $m \times n$  matrix, then letting  $\mathbf{e}_i \in \mathbb{R}^n$  and  $\mathbf{e}_j \in \mathbb{R}^n$ ,

$$\mathbf{e}_i^T A \mathbf{e}_j = A_{ij} \quad (19.23)$$

**Proof:** First note that  $\mathbf{e}_i^T$  is a  $1 \times n$  matrix and  $\mathbf{v}$  is an  $n \times 1$  matrix so the above multiplication in (19.22) makes perfect sense. It equals

$$(0, \dots, 1, \dots, 0) \begin{pmatrix} v_1 \\ \vdots \\ v_i \\ \vdots \\ v_n \end{pmatrix} = v_i$$

as claimed.

Consider (19.23). From the definition of matrix multiplication using the repeated index summation convention, and noting that  $(\mathbf{e}_j)_k = \delta_{kj}$

$$\mathbf{e}_i^T A \mathbf{e}_j = \mathbf{e}_i^T \begin{pmatrix} A_{1k} (\mathbf{e}_j)_k \\ \vdots \\ A_{ik} (\mathbf{e}_j)_k \\ \vdots \\ A_{mk} (\mathbf{e}_j)_k \end{pmatrix} = \mathbf{e}_i^T \begin{pmatrix} A_{1j} \\ \vdots \\ A_{ij} \\ \vdots \\ A_{mj} \end{pmatrix} = A_{ij}$$

by the first part of the lemma. This proves the lemma.

**Theorem 19.3.4** *Let  $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a linear transformation. Then there exists a unique  $m \times n$  matrix,  $A$  such that*

$$A\mathbf{x} = L\mathbf{x}$$

for all  $\mathbf{x} \in \mathbb{R}^n$ . The  $ik^{th}$  entry of this matrix is given by

$$\mathbf{e}_i^T L \mathbf{e}_k \quad (19.24)$$

**Proof:** By the lemma,

$$(L\mathbf{x})_i = \mathbf{e}_i^T L\mathbf{x} = \mathbf{e}_i^T x_k L \mathbf{e}_k = (\mathbf{e}_i^T L \mathbf{e}_k) x_k.$$

Let  $A_{ik} = \mathbf{e}_i^T L \mathbf{e}_k$ , to prove the existence part of the theorem.

To verify uniqueness, suppose  $B\mathbf{x} = A\mathbf{x} = L\mathbf{x}$  for all  $\mathbf{x} \in \mathbb{R}^n$ . Then in particular, this is true for  $\mathbf{x} = \mathbf{e}_j$  and then multiply on the left by  $\mathbf{e}_i^T$  to obtain

$$B_{ij} = \mathbf{e}_i^T B \mathbf{e}_j = \mathbf{e}_i^T A \mathbf{e}_j = A_{ij}$$

showing  $A = B$ . This proves uniqueness.

**Corollary 19.3.5** *A linear transformation,  $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is completely determined by the vectors  $\{L\mathbf{e}_1, \dots, L\mathbf{e}_n\}$ .*

**Proof:** This follows immediately from the above theorem. The unique matrix determining the linear transformation which is given in (19.24) depends only on these vectors.

This theorem shows that any linear transformation defined on  $\mathbb{R}^n$  can always be considered as a matrix. Therefore, the terms “linear transformation” and “matrix” will be used interchangeably. For example, to say a matrix is one to one, means the linear transformation determined by the matrix is one to one.

**Example 19.3.6** *Find the linear transformation,  $L : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  which has the property that  $L\mathbf{e}_1 = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$  and  $L\mathbf{e}_2 = \begin{pmatrix} 1 \\ 3 \end{pmatrix}$ . From the above theorem and corollary, this linear transformation is that determined by matrix multiplication by the matrix*

$$\begin{pmatrix} 2 & 1 \\ 1 & 3 \end{pmatrix}.$$

### 19.3.1 Least Squares Problems

**Definition 19.3.7** Let  $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$  be a set of vectors in  $\mathbb{R}^n$ . Define

$$\text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_k\} \equiv \left\{ \sum_{i=1}^k c_i \mathbf{x}_i : c_1, \dots, c_k \in \mathbb{R} \right\}$$

In words  $\text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$  consists of all vectors which can be obtained by adding up scalars times the vectors in the set  $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ . These sums are called linear combinations.

The following lemma is very significant, being a special case of something called the Gram Schmidt procedure. Before reading further, you should review the dot product and its properties on Page 311. Also recall that

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}.$$

**Lemma 19.3.8** Let  $A$  be an  $m \times n$  matrix,  $A \neq 0$ . Then there exist vectors,  $\{\mathbf{f}_1, \dots, \mathbf{f}_r\}$  such that  $\mathbf{f}_i \cdot \mathbf{f}_j = \delta_{ij}$  and  $\text{span}\{\mathbf{f}_1, \dots, \mathbf{f}_r\} = A(\mathbb{R}^n) \equiv \{A\mathbf{x} : \mathbf{x} \in \mathbb{R}^n\}$ .

**Proof:** Let

$$\mathbf{f}_1 = \frac{A\mathbf{e}_{i_1}}{|A\mathbf{e}_{i_1}|}$$

where  $\mathbf{e}_{i_1}$  is the first vector in the list,  $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  with the property that  $A\mathbf{e}_i \neq \mathbf{0}$ . Thus

$$\text{span}\{\mathbf{f}_1\} = \text{span}\{A\mathbf{e}_{i_1}\} = \text{span}\{A\mathbf{e}_1, \dots, A\mathbf{e}_{i_1}\}$$

because  $A\mathbf{e}_k = \mathbf{0}$  for  $k < i_1$ . Now suppose  $\{\mathbf{f}_1, \dots, \mathbf{f}_k\}$  have been chosen in such a way that  $\mathbf{f}_i \cdot \mathbf{f}_j = \delta_{ij}$  and

$$\text{span}\{A\mathbf{e}_{i_1}, \dots, A\mathbf{e}_{i_k}\} = \text{span}\{A\mathbf{e}_1, A\mathbf{e}_2, \dots, A\mathbf{e}_{i_k}\} = \text{span}\{\mathbf{f}_1, \dots, \mathbf{f}_k\}.$$

This was just accomplished for  $k = 1$ . There are two cases to consider.

First it could happen that  $A\mathbf{e}_j \in \text{span}\{\mathbf{f}_1, \dots, \mathbf{f}_k\}$  for all  $j$ . In this case, stop and consider  $\text{span}\{\mathbf{f}_1, \dots, \mathbf{f}_k\}$ . For  $\mathbf{x} \in \mathbb{R}^n$  arbitrary,

$$\begin{aligned} A\mathbf{x} &= A \left( \sum_{i=1}^n x_i \mathbf{e}_i \right) = \sum_{i=1}^k x_i A\mathbf{e}_i \\ &= \sum_{i=1}^n x_i \overbrace{\sum_{j=1}^k y_{ij} \mathbf{f}_j}^{=A\mathbf{e}_i} = \sum_{j=1}^k \left( \sum_{i=1}^n x_i y_{ij} \right) \mathbf{f}_j \\ &\in \text{span}\{\mathbf{f}_1, \dots, \mathbf{f}_k\} \end{aligned}$$

showing that  $\text{span}\{\mathbf{f}_1, \dots, \mathbf{f}_k\} = A(\mathbb{R}^n)$ .

The other case is that  $A\mathbf{e}_j \notin \text{span}\{\mathbf{f}_1, \dots, \mathbf{f}_k\}$  for some  $j$ . (By construction,  $A\mathbf{e}_j$  is in  $\text{span}\{\mathbf{f}_1, \dots, \mathbf{f}_k\}$  for all  $j \leq i_k$ .) Let  $i_{k+1}$  be the smallest index larger than  $i_k$  such that  $A\mathbf{e}_{i_{k+1}} \notin \text{span}\{\mathbf{f}_1, \dots, \mathbf{f}_k\}$ . Then define

$$\mathbf{f}_{k+1} \equiv \frac{A\mathbf{e}_{i_{k+1}} - \sum_{j=1}^k (A\mathbf{e}_{i_{k+1}} \cdot \mathbf{f}_j) \mathbf{f}_j}{\left| A\mathbf{e}_{i_{k+1}} - \sum_{j=1}^k (A\mathbf{e}_{i_{k+1}} \cdot \mathbf{f}_j) \mathbf{f}_j \right|}. \quad (19.25)$$

Note the denominator is non zero because of the assumption that  $A\mathbf{e}_{i_{k+1}} \notin \text{span}\{\mathbf{f}_1, \dots, \mathbf{f}_k\}$ . Then

$$\text{span}\{\mathbf{f}_1, \dots, \mathbf{f}_{k+1}\} = \text{span}\{A\mathbf{e}_{i_1}, \dots, A\mathbf{e}_{i_{k+1}}\} = \text{span}\{A\mathbf{e}_1, \dots, A\mathbf{e}_{i_{k+1}}\} \quad (19.26)$$

Also, from the properties of the dot product on Page 311, for  $l < k + 1$ ,

$$\begin{aligned} \mathbf{f}_{k+1} \cdot \mathbf{f}_l &= \frac{A\mathbf{e}_{i_{k+1}} \cdot \mathbf{f}_l - \sum_{j=1}^k (A\mathbf{e}_{i_{k+1}} \cdot \mathbf{f}_j) \overbrace{\mathbf{f}_j \cdot \mathbf{f}_l}^{\delta_{jl}}}{\left| A\mathbf{e}_{i_{k+1}} - \sum_{j=1}^k (A\mathbf{e}_{i_{k+1}} \cdot \mathbf{f}_j) \mathbf{f}_j \right|}} \\ &= \frac{A\mathbf{e}_{i_{k+1}} \cdot \mathbf{f}_l - A\mathbf{e}_{i_{k+1}} \cdot \mathbf{f}_l}{\left| A\mathbf{e}_{i_{k+1}} - \sum_{j=1}^k (A\mathbf{e}_{i_{k+1}} \cdot \mathbf{f}_j) \mathbf{f}_j \right|} = 0 \end{aligned}$$

Therefore, since  $|\mathbf{f}_{k+1}| = 1$ , it follows that  $\{\mathbf{f}_1, \dots, \mathbf{f}_{k+1}\}$  has the property that  $\mathbf{f}_i \cdot \mathbf{f}_j = \delta_{ij}$  and each  $\mathbf{f}_j \in A(\mathbb{R}^n)$ . Continue until the first case is obtained. This must happen in no more than  $n$  applications of the above process because for each  $j$ ,  $i_j < i_{j+1}$  and  $i_j$  can't be any larger than  $n$ . This proves the lemma.

The following is a fundamental existence theorem which is based on the above lemma.

**Theorem 19.3.9** *Let  $\mathbf{y} \in \mathbb{R}^m$  and let  $A$  be an  $m \times n$  matrix. Then there exists  $\mathbf{x} \in \mathbb{R}^n$  minimizing the function,  $|\mathbf{y} - A\mathbf{x}|^2$ . Furthermore,  $\mathbf{x}$  minimizes this function if and only if*

$$(\mathbf{y} - A\mathbf{x}) \cdot A\mathbf{w} = 0$$

for all  $\mathbf{w} \in \mathbb{R}^n$ .

**Proof:** Let  $\{\mathbf{f}_1, \dots, \mathbf{f}_r\}$  be the vectors of Lemma 19.3.8 with  $\text{span}\{\mathbf{f}_1, \dots, \mathbf{f}_r\} = A(\mathbb{R}^n)$ . Since  $A(\mathbb{R}^n) = \text{span}\{\mathbf{f}_1, \dots, \mathbf{f}_r\}$ , it follows that if you can find  $y_1, \dots, y_r$  in such a way as to minimize

$$\left| \mathbf{y} - \sum_{k=1}^r y_k \mathbf{f}_k \right|^2,$$

then letting  $A\mathbf{x} = \sum_{k=1}^r y_k \mathbf{f}_k$ , it will follow that this  $\mathbf{x}$  is the desired solution. Let  $y_1, \dots, y_r$  be a list of scalars. Then from the definition of  $|\cdot|$  and the properties of the dot product,

$$\begin{aligned} \left| \mathbf{y} - \sum_{k=1}^r y_k \mathbf{f}_k \right|^2 &= \left( \mathbf{y} - \sum_{k=1}^r y_k \mathbf{f}_k \right) \cdot \left( \mathbf{y} - \sum_{k=1}^r y_k \mathbf{f}_k \right) \\ &= |\mathbf{y}|^2 - 2 \sum_{k=1}^r y_k (\mathbf{y} \cdot \mathbf{f}_k) + \sum_k \sum_l y_k y_l \overbrace{\mathbf{f}_k \cdot \mathbf{f}_l}^{\delta_{kl}} \\ &= |\mathbf{y}|^2 - 2 \sum_{k=1}^r y_k (\mathbf{y} \cdot \mathbf{f}_k) + \sum_{k=1}^r y_k^2 \\ &= |\mathbf{y}|^2 + \sum_{k=1}^r y_k^2 - 2y_k (\mathbf{y} \cdot \mathbf{f}_k) \end{aligned}$$

Now complete the square to obtain

$$\begin{aligned} &= |\mathbf{y}|^2 + \sum_{k=1}^r \left( y_k^2 - 2y_k (\mathbf{y} \cdot \mathbf{f}_k) + (\mathbf{y} \cdot \mathbf{f}_k)^2 \right) - \sum_{k=1}^r (\mathbf{y} \cdot \mathbf{f}_k)^2 \\ &= |\mathbf{y}|^2 + \sum_{k=1}^r (y_k - (\mathbf{y} \cdot \mathbf{f}_k))^2 - \sum_{k=1}^r (\mathbf{y} \cdot \mathbf{f}_k)^2. \end{aligned}$$



This shows that the minimum is obtained when  $y_k = (\mathbf{y} \cdot \mathbf{f}_k)$ . This proves the existence part of the Theorem.

To verify the other part, let  $t \in \mathbb{R}$  and consider

$$\begin{aligned} |\mathbf{y} - A(\mathbf{x} + t\mathbf{w})|^2 &= (\mathbf{y} - A\mathbf{x} - tA\mathbf{w}) \cdot (\mathbf{y} - A\mathbf{x} - tA\mathbf{w}) \\ &= |\mathbf{y} - A\mathbf{x}|^2 - 2t(\mathbf{y} - A\mathbf{x}) \cdot A\mathbf{w} + t^2 |A\mathbf{w}|^2. \end{aligned}$$

Then from the above equation,  $|\mathbf{y} - A\mathbf{x}|^2 \leq |\mathbf{y} - A\mathbf{z}|^2$  for all  $\mathbf{z} \in \mathbb{R}^n$  if and only if for all  $\mathbf{w} \in \mathbb{R}^n$  and  $t \in \mathbb{R}$

$$|\mathbf{y} - A\mathbf{x}|^2 - 2t(\mathbf{y} - A\mathbf{x}) \cdot A\mathbf{w} + t^2 |\mathbf{w}|^2 \geq |\mathbf{y} - A\mathbf{x}|^2$$

and this happens if and only if for all  $t \in \mathbb{R}$  and  $\mathbf{w} \in \mathbb{R}^n$ ,

$$-2t(\mathbf{y} - A\mathbf{x}) \cdot A\mathbf{w} + t^2 |A\mathbf{w}|^2 \geq 0,$$

which occurs if and only if  $(\mathbf{y} - A\mathbf{x}) \cdot A\mathbf{w} = 0$  for all  $\mathbf{w} \in \mathbb{R}^n$ . (Why?) This proves the theorem.

**Lemma 19.3.10** *Let  $A$  be an  $m \times n$  matrix. Then*

$$A\mathbf{x} \cdot \mathbf{y} = \mathbf{x} \cdot A^T \mathbf{y}$$

**Proof:** This follows from the definition. Using the repeated index summation convention,

$$\begin{aligned} A\mathbf{x} \cdot \mathbf{y} &= A_{ij} x_j y_i \\ &= x_j A_{ji}^T y_i \\ &= \mathbf{x} \cdot A^T \mathbf{y}. \end{aligned}$$

### 19.3.2 The Least Squares Regression Line

As an important application of the above theorem is the problem of finding the least squares regression line in statistics. Suppose you are given points in the plane,  $\{(x_i, y_i)\}_{i=1}^n$  and you would like to find constants  $m$  and  $b$  such that the line  $y = mx + b$  goes through all these points. Of course this will be impossible in general. Therefore, try to find  $m, b$  to get as close as possible. The desired system is

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_1 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix} \begin{pmatrix} m \\ b \end{pmatrix}$$

which is of the form  $\mathbf{y} = A\mathbf{x}$  and it is desired to choose  $m$  and  $b$  to make

$$\left| A \begin{pmatrix} m \\ b \end{pmatrix} - \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \right|^2$$

as small as possible. According to Theorem 19.3.9, the best values for  $m$  and  $b$  occur as the solution to

$$A^T A \begin{pmatrix} m \\ b \end{pmatrix} = A^T \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

where

$$A = \begin{pmatrix} x_1 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix}.$$

Thus, computing  $A^T A$ ,

$$\begin{pmatrix} \sum_{i=1}^n x_i^2 & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & n \end{pmatrix} \begin{pmatrix} m \\ b \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n x_i y_i \\ \sum_{i=1}^n y_i \end{pmatrix}$$

Solving this system of equations for  $m$  and  $b$ ,

$$m = \frac{-(\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i) + (\sum_{i=1}^n x_i y_i) n}{(\sum_{i=1}^n x_i^2) n - (\sum_{i=1}^n x_i)^2}$$

and

$$b = \frac{-(\sum_{i=1}^n x_i) \sum_{i=1}^n x_i y_i + (\sum_{i=1}^n y_i) \sum_{i=1}^n x_i^2}{(\sum_{i=1}^n x_i^2) n - (\sum_{i=1}^n x_i)^2}.$$

One could clearly do a least squares fit for curves of the form  $y = ax^2 + bx + c$  in the same way. In this case you want to solve as well as possible for  $a, b$ , and  $c$  the system

$$\begin{pmatrix} x_1^2 & x_1 & 1 \\ \vdots & \vdots & \vdots \\ x_n^2 & x_n & 1 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

and one would use the same technique as above. Many other similar problems are important, including many in higher dimensions and they are all solved the same way.

### 19.3.3 The Fredholm Alternative

The next major result is called the Fredholm alternative. It comes from Theorem 19.3.9 and Lemma 19.3.10.

**Theorem 19.3.11** *Let  $A$  be an  $m \times n$  matrix. Then there exists  $\mathbf{x} \in \mathbb{R}^n$  such that  $A\mathbf{x} = \mathbf{y}$  if and only if whenever  $A^T \mathbf{z} = \mathbf{0}$  it follows that  $\mathbf{z} \cdot \mathbf{y} = 0$ .*

**Proof:** First suppose that for some  $\mathbf{x} \in \mathbb{R}^n$ ,  $A\mathbf{x} = \mathbf{y}$ . Then letting  $A^T \mathbf{z} = \mathbf{0}$  and using the above lemma,

$$\mathbf{y} \cdot \mathbf{z} = A\mathbf{x} \cdot \mathbf{z} = \mathbf{x} \cdot A^T \mathbf{z} = \mathbf{x} \cdot \mathbf{0} = 0.$$

This proves half the theorem.

To do the other half, suppose that whenever,  $A^T \mathbf{z} = \mathbf{0}$  it follows that  $\mathbf{z} \cdot \mathbf{y} = 0$ . It is necessary to show there exists  $\mathbf{x} \in \mathbb{R}^n$  such that  $\mathbf{y} = A\mathbf{x}$ . From Theorem 19.3.9 there exists  $\mathbf{x}$  minimizing  $|\mathbf{y} - A\mathbf{x}|^2$  which therefore satisfies

$$(\mathbf{y} - A\mathbf{x}) \cdot A\mathbf{w} = 0 \tag{19.27}$$

for all  $\mathbf{w} \in \mathbb{R}^n$ . Therefore, for all  $\mathbf{w} \in \mathbb{R}^n$ ,

$$A^T (\mathbf{y} - A\mathbf{x}) \cdot \mathbf{w} = 0$$

which shows that  $A^T (\mathbf{y} - A\mathbf{x}) = \mathbf{0}$ . (Why?) Therefore, by assumption,

$$(\mathbf{y} - A\mathbf{x}) \cdot \mathbf{y} = 0.$$

Now by (19.27) with  $\mathbf{w} = \mathbf{x}$ ,

$$(\mathbf{y} - A\mathbf{x}) \cdot (\mathbf{y} - A\mathbf{x}) = (\mathbf{y} - A\mathbf{x}) \cdot \mathbf{y} - (\mathbf{y} - A\mathbf{x}) \cdot A\mathbf{x} = 0$$

showing that  $\mathbf{y} = A\mathbf{x}$ . This proves the theorem.

The following corollary is also called the Fredholm alternative.

**Corollary 19.3.12** *Let  $A$  be an  $m \times n$  matrix. Then  $A$  is onto if and only if  $A^T$  is one to one.*

**Proof:** Suppose first  $A$  is onto. Then by Theorem 19.3.11, it follows that for all  $\mathbf{y} \in \mathbb{R}^m$ ,  $\mathbf{y} \cdot \mathbf{z} = 0$  whenever  $A^T \mathbf{z} = \mathbf{0}$ . Therefore, let  $\mathbf{y} = \mathbf{z}$  where  $A^T \mathbf{z} = \mathbf{0}$  and conclude that  $\mathbf{z} \cdot \mathbf{z} = 0$  whenever  $A^T \mathbf{z} = \mathbf{0}$ . If  $A^T \mathbf{x} = A^T \mathbf{y}$ , then  $A^T (\mathbf{x} - \mathbf{y}) = \mathbf{0}$  and so  $\mathbf{x} - \mathbf{y} = \mathbf{0}$ . Thus  $A^T$  is one to one.

Now let  $\mathbf{y} \in \mathbb{R}^m$  be given.  $\mathbf{y} \cdot \mathbf{z} = 0$  whenever  $A^T \mathbf{z} = \mathbf{0}$  because, since  $A^T$  is assumed to be one to one, and  $\mathbf{0}$  is a solution to this equation, it must be the only solution. Therefore, by Theorem 19.3.11 there exists  $\mathbf{x}$  such that  $A\mathbf{x} = \mathbf{y}$  therefore,  $A$  is onto.

## 19.4 Exercises

1. The proof of Theorem 19.3.9 concluded with the following observation. If  $-ta + t^2b \geq 0$  for all  $t \in \mathbb{R}$  and  $b \geq 0$ , then  $a = 0$ . Why is this so?
2. In the proof of Lemma 19.3.8 explain the assertion (19.26).
3. In the proof of Theorem 19.3.11 the following argument was used. If  $\mathbf{x} \cdot \mathbf{w} = 0$  for all  $\mathbf{w} \in \mathbb{R}^n$ , then  $\mathbf{x} = \mathbf{0}$ . Why is this so?
4. Suppose  $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a linear transformation. Show the following are equivalent.
  - (a)  $L\mathbf{x} = \mathbf{0}$  implies  $\mathbf{x} = \mathbf{0}$ .
  - (b)  $L$  is one to one.
5. Using Corollary 19.3.12 and Problem 4, show that an  $m \times n$  matrix is onto if and only if its transpose is one to one.
6. Suppose  $A$  is a  $3 \times 2$  matrix. Is it possible that  $A^T$  is one to one? What does this say about  $A$  being onto? Prove your answer.
7. Explain why there always exists a solution to the equation,  $A^T \mathbf{y} = A^T A \mathbf{x}$  and also explain why this is called a least squares solution to the equation,  $\mathbf{y} = A\mathbf{x}$ .
8. Referring to Problem 7, find the least squares solution to the following system.

$$\begin{aligned} x + 2y &= 1 \\ 2x + 3y &= 2 \\ 3x + 5y &= 4 \end{aligned}$$

9. You are doing experiments and have obtained the ordered pairs,  $(0, 1)$ ,  $(1, 2)$ ,  $(2, 3.5)$ , and  $(3, 4)$ . Find  $m$  and  $b$  such that  $y = mx + b$  approximates these four points as well as possible. Now do the same thing for  $y = ax^2 + bx + c$ , finding  $a, b$ , and  $c$  to give the best approximation.

10. Suppose you have several ordered triples,  $(x_i, y_i, z_i)$ . Describe how to find a polynomial,

$$z = a + bx + cy + dxy + ex^2 + fy^2$$

for example giving the best fit to the given ordered triples. Is there any reason you have to use a polynomial? Would similar approaches work for other combinations of functions just as well?

## 19.5 Moving Coordinate Systems

The study of moving coordinate systems gives a non trivial example of the usefulness of the above ideas involving linear transformations and matrices. To begin with, here is the concept of the product rule extended to matrix multiplication.

**Definition 19.5.1** Let  $A(t)$  be an  $m \times n$  matrix. Say  $A(t) = (A_{ij}(t))$ . Suppose also that  $A_{ij}(t)$  is a differentiable function for all  $i, j$ . Then define  $A'(t) \equiv (A'_{ij}(t))$ . That is,  $A'(t)$  is the matrix which consists of replacing each entry by its derivative. Such an  $m \times n$  matrix in which the entries are differentiable functions is called a differentiable matrix.

The next lemma is just a version of the product rule.

**Lemma 19.5.2** Let  $A(t)$  be an  $m \times n$  matrix and let  $B(t)$  be an  $n \times p$  matrix with the property that all the entries of these matrices are differentiable functions. Then

$$(A(t)B(t))' = A'(t)B(t) + A(t)B'(t).$$

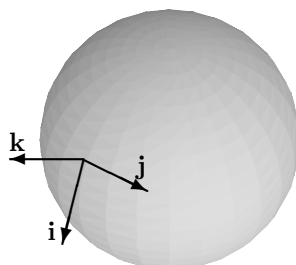
**Proof:**  $(A(t)B(t))' = (C'_{ij}(t))$  where  $C_{ij}(t) = A_{ik}(t)B_{kj}(t)$  and the repeated index summation convention is being used. Therefore,

$$\begin{aligned} C'_{ij}(t) &= A'_{ik}(t)B_{kj}(t) + A_{ik}(t)B'_{kj}(t) \\ &= (A'(t)B(t))_{ij} + (A(t)B'(t))_{ij} \\ &= (A'(t)B(t) + A(t)B'(t))_{ij} \end{aligned}$$

Therefore, the  $ij^{th}$  entry of  $A(t)B(t)$  equals the  $ij^{th}$  entry of  $A'(t)B(t) + A(t)B'(t)$  and this proves the lemma.

### 19.5.1 The Coriolis Acceleration

Imagine a point on the surface of the earth. Now consider unit vectors, one pointing South, one pointing East and one pointing directly away from the center of the earth.



Denote the first as  $\mathbf{i}$ , the second as  $\mathbf{j}$  and the third as  $\mathbf{k}$ . If you are standing on the earth you will consider these vectors as fixed, but of course they are not. As the earth turns, they change direction and so each is in reality a function of  $t$ . Nevertheless, it is with respect to these apparently fixed vectors that you wish to understand acceleration, velocities, and displacements.

In general, let  $\mathbf{i}^*, \mathbf{j}^*, \mathbf{k}^*$  be the usual fixed vectors in space and let  $\mathbf{i}(t), \mathbf{j}(t), \mathbf{k}(t)$  be an orthonormal basis of vectors for each  $t$ , like the vectors described in the first paragraph. It is assumed these vectors are  $C^1$  functions of  $t$ . Letting the positive  $x$  axis extend in the direction of  $\mathbf{i}(t)$ , the positive  $y$  axis extend in the direction of  $\mathbf{j}(t)$ , and the positive  $z$  axis extend in the direction of  $\mathbf{k}(t)$ , yields a moving coordinate system. Now let  $\mathbf{u}$  be a vector and let  $t_0$  be some reference time. For example you could let  $t_0 = 0$ . Then define the components of  $\mathbf{u}$  with respect to these vectors,  $\mathbf{i}, \mathbf{j}, \mathbf{k}$  at time  $t_0$  as

$$\mathbf{u} \equiv u^1 \mathbf{i}(t_0) + u^2 \mathbf{j}(t_0) + u^3 \mathbf{k}(t_0).$$

Let  $\mathbf{u}(t)$  be defined as the vector which has the same components with respect to  $\mathbf{i}, \mathbf{j}, \mathbf{k}$  but at time  $t$ . Thus

$$\mathbf{u}(t) \equiv u^1 \mathbf{i}(t) + u^2 \mathbf{j}(t) + u^3 \mathbf{k}(t).$$

and the vector has changed although the components have not.

This is exactly the situation in the case of the apparently fixed basis vectors on the earth if  $\mathbf{u}$  is a position vector from the given spot on the earth's surface to a point regarded as fixed with the earth due to its keeping the same coordinates relative to the coordinate axes which are fixed with the earth. Now define a linear transformation  $Q(t)$  mapping  $\mathbb{R}^3$  to  $\mathbb{R}^3$  by

$$Q(t) \mathbf{u} \equiv u^1 \mathbf{i}(t) + u^2 \mathbf{j}(t) + u^3 \mathbf{k}(t)$$

where

$$\mathbf{u} \equiv u^1 \mathbf{i}(t_0) + u^2 \mathbf{j}(t_0) + u^3 \mathbf{k}(t_0)$$

Thus letting  $\mathbf{v}$  be a vector defined in the same manner as  $\mathbf{u}$  and  $\alpha, \beta$ , scalars,

$$\begin{aligned} Q(t)(\alpha \mathbf{u} + \beta \mathbf{v}) &\equiv (\alpha u^1 + \beta v^1) \mathbf{i}(t) + (\alpha u^2 + \beta v^2) \mathbf{j}(t) + (\alpha u^3 + \beta v^3) \mathbf{k}(t) \\ &= (\alpha u^1 \mathbf{i}(t) + \alpha u^2 \mathbf{j}(t) + \alpha u^3 \mathbf{k}(t)) + (\beta v^1 \mathbf{i}(t) + \beta v^2 \mathbf{j}(t) + \beta v^3 \mathbf{k}(t)) \\ &= \alpha (u^1 \mathbf{i}(t) + u^2 \mathbf{j}(t) + u^3 \mathbf{k}(t)) + \beta (v^1 \mathbf{i}(t) + v^2 \mathbf{j}(t) + v^3 \mathbf{k}(t)) \\ &\equiv \alpha Q(t) \mathbf{u} + \beta Q(t) \mathbf{v} \end{aligned}$$

showing that  $Q(t)$  is a linear transformation. Also,  $Q(t)$  preserves all distances because, since the vectors,  $\mathbf{i}(t), \mathbf{j}(t), \mathbf{k}(t)$  form an orthonormal set,

$$|Q(t) \mathbf{u}| = \left( \sum_{i=1}^3 (u^i)^2 \right)^{1/2} = |\mathbf{u}|.$$

**Lemma 19.5.3** Suppose  $Q(t)$  is a real, differentiable  $n \times n$  matrix which preserves distances. Then  $Q(t) Q(t)^T = Q(t)^T Q(t) = I$ . Also, if  $\mathbf{u}(t) \equiv Q(t) \mathbf{u}$ , then there exists a vector,  $\boldsymbol{\Omega}(t)$  such that

$$\mathbf{u}'(t) = \boldsymbol{\Omega}(t) \times \mathbf{u}(t).$$

**Proof:** Recall that  $(\mathbf{z} \cdot \mathbf{w}) = \frac{1}{4} (|\mathbf{z} + \mathbf{w}|^2 - |\mathbf{z} - \mathbf{w}|^2)$ . Therefore,

$$\begin{aligned} (Q(t) \mathbf{u} \cdot Q(t) \mathbf{w}) &= \frac{1}{4} (|Q(t) (\mathbf{u} + \mathbf{w})|^2 - |Q(t) (\mathbf{u} - \mathbf{w})|^2) \\ &= \frac{1}{4} (|\mathbf{u} + \mathbf{w}|^2 - |\mathbf{u} - \mathbf{w}|^2) \\ &= (\mathbf{u} \cdot \mathbf{w}). \end{aligned}$$

This implies

$$(Q(t)^T Q(t) \mathbf{u} \cdot \mathbf{w}) = (\mathbf{u} \cdot \mathbf{w})$$

for all  $\mathbf{u}, \mathbf{w}$ . Therefore,  $Q(t)^T Q(t) \mathbf{u} = \mathbf{u}$  and so  $Q(t)^T Q(t) = Q(t) Q(t)^T = I$ . This proves the first part of the lemma.

It follows from the product rule, Lemma 19.5.2 that

$$Q'(t) Q(t)^T + Q(t) Q'(t)^T = 0$$

and so

$$Q'(t) Q(t)^T = - (Q'(t) Q(t)^T)^T. \quad (19.28)$$

From the definition,  $Q(t) \mathbf{u} = \mathbf{u}(t)$ ,

$$\mathbf{u}'(t) = Q'(t) \mathbf{u} = Q'(t) \overbrace{Q(t)^T}^{=\mathbf{u}} \mathbf{u}(t).$$

Then writing the matrix of  $Q'(t) Q(t)^T$  with respect to fixed in space orthonormal basis vectors,  $\mathbf{i}^*, \mathbf{j}^*, \mathbf{k}^*$ , where these are the usual basis vectors for  $\mathbb{R}^3$ , it follows from (19.28) that the matrix of  $Q'(t) Q(t)^T$  is of the form

$$\begin{pmatrix} 0 & -\omega_3(t) & \omega_2(t) \\ \omega_3(t) & 0 & -\omega_1(t) \\ -\omega_2(t) & \omega_1(t) & 0 \end{pmatrix}$$

for some time dependent scalars,  $\omega_i$ . Therefore,

$$\begin{aligned} \begin{pmatrix} u^1 \\ u^2 \\ u^3 \end{pmatrix}'(t) &= \begin{pmatrix} 0 & -\omega_3(t) & \omega_2(t) \\ \omega_3(t) & 0 & -\omega_1(t) \\ -\omega_2(t) & \omega_1(t) & 0 \end{pmatrix} \begin{pmatrix} u^1 \\ u^2 \\ u^3 \end{pmatrix}(t) \\ &= \begin{pmatrix} \omega_2(t) u^3(t) - \omega_3(t) u^2(t) \\ \omega_3(t) u^1(t) - \omega_1(t) u^3(t) \\ \omega_1(t) u^2(t) - \omega_2(t) u^1(t) \end{pmatrix} \end{aligned}$$

where the  $u^i$  are the components of the vector  $\mathbf{u}(t)$  in terms of the fixed vectors  $\mathbf{i}^*, \mathbf{j}^*, \mathbf{k}^*$ . Therefore,

$$\mathbf{u}'(t) = \boldsymbol{\Omega}(t) \times \mathbf{u}(t) = Q'(t) Q(t)^T \mathbf{u}(t) \quad (19.29)$$

where

$$\boldsymbol{\Omega}(t) = \omega_1(t) \mathbf{i}^* + \omega_2(t) \mathbf{j}^* + \omega_3(t) \mathbf{k}^*.$$

because

$$\boldsymbol{\Omega}(t) \times \mathbf{u}(t) \equiv \begin{vmatrix} \mathbf{i}^* & \mathbf{j}^* & \mathbf{k}^* \\ \omega_1 & \omega_2 & \omega_3 \\ u^1 & u^2 & u^3 \end{vmatrix} \equiv$$

$$\mathbf{i}^* (\omega_2 u^3 - \omega_3 u^2) + \mathbf{j}^* (\omega_3 u^1 - \omega_1 u^3) + \mathbf{k}^* (\omega_1 u^2 - \omega_2 u^1).$$

This proves the lemma and yields the existence part of the following theorem.

**Theorem 19.5.4** Let  $\mathbf{i}(t), \mathbf{j}(t), \mathbf{k}(t)$  be as described. Then there exists a unique vector  $\boldsymbol{\Omega}(t)$  such that if  $\mathbf{u}(t)$  is a vector whose components are constant with respect to  $\mathbf{i}(t), \mathbf{j}(t), \mathbf{k}(t)$ , then

$$\mathbf{u}'(t) = \boldsymbol{\Omega}(t) \times \mathbf{u}(t).$$

**Proof:** It only remains to prove uniqueness. Suppose  $\boldsymbol{\Omega}_1$  also works. Then  $\mathbf{u}(t) = Q(t) \mathbf{u}$  and so  $\mathbf{u}'(t) = Q'(t) \mathbf{u}$  and

$$Q'(t) \mathbf{u} = \boldsymbol{\Omega} \times Q(t) \mathbf{u} = \boldsymbol{\Omega}_1 \times Q(t) \mathbf{u}$$

for all  $\mathbf{u}$ . Therefore,

$$(\boldsymbol{\Omega} - \boldsymbol{\Omega}_1) \times Q(t) \mathbf{u} = \mathbf{0}$$

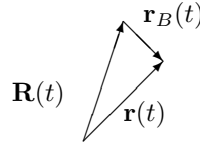
for all  $\mathbf{u}$  and since  $Q(t)$  is one to one and onto, this implies  $(\boldsymbol{\Omega} - \boldsymbol{\Omega}_1) \times \mathbf{w} = \mathbf{0}$  for all  $\mathbf{w}$  and thus  $\boldsymbol{\Omega} - \boldsymbol{\Omega}_1 = \mathbf{0}$ . This proves the theorem.

Now let  $\mathbf{R}(t)$  be a position vector and let

$$\mathbf{r}(t) = \mathbf{R}(t) + \mathbf{r}_B(t)$$

where

$$\mathbf{r}_B(t) \equiv x(t)\mathbf{i}(t) + y(t)\mathbf{j}(t) + z(t)\mathbf{k}(t).$$



In the example of the earth,  $\mathbf{R}(t)$  is the position vector of a point  $\mathbf{p}(t)$  on the earth's surface and  $\mathbf{r}_B(t)$  is the position vector of another point from  $\mathbf{p}(t)$ , thus regarding  $\mathbf{p}(t)$  as the origin.  $\mathbf{r}_B(t)$  is the position vector of a point as perceived by the observer on the earth with respect to the vectors he thinks of as fixed. Similarly,  $\mathbf{v}_B(t)$  and  $\mathbf{a}_B(t)$  will be the velocity and acceleration relative to  $\mathbf{i}(t), \mathbf{j}(t), \mathbf{k}(t)$ , and so  $\mathbf{v}_B = x'\mathbf{i} + y'\mathbf{j} + z'\mathbf{k}$  and  $\mathbf{a}_B = x''\mathbf{i} + y''\mathbf{j} + z''\mathbf{k}$ . Then

$$\mathbf{v} \equiv \mathbf{r}' = \mathbf{R}' + x'\mathbf{i} + y'\mathbf{j} + z'\mathbf{k} + x\mathbf{i}' + y\mathbf{j}' + z\mathbf{k}'.$$

By (19.29), if  $\mathbf{e} \in \{\mathbf{i}, \mathbf{j}, \mathbf{k}\}$ ,  $\mathbf{e}' = \boldsymbol{\Omega} \times \mathbf{e}$  because the components of these vectors with respect to  $\mathbf{i}, \mathbf{j}, \mathbf{k}$  are constant. Therefore,

$$\begin{aligned} x\mathbf{i}' + y\mathbf{j}' + z\mathbf{k}' &= x\boldsymbol{\Omega} \times \mathbf{i} + y\boldsymbol{\Omega} \times \mathbf{j} + z\boldsymbol{\Omega} \times \mathbf{k} \\ &= \boldsymbol{\Omega} \times (x\mathbf{i} + y\mathbf{j} + z\mathbf{k}) \end{aligned}$$

and consequently,

$$\mathbf{v} = \mathbf{R}' + x'\mathbf{i} + y'\mathbf{j} + z'\mathbf{k} + \boldsymbol{\Omega} \times \mathbf{r}_B = \mathbf{R}' + x'\mathbf{i} + y'\mathbf{j} + z'\mathbf{k} + \boldsymbol{\Omega} \times (x\mathbf{i} + y\mathbf{j} + z\mathbf{k}).$$

Now consider the acceleration. Quantities which are relative to the moving coordinate system and quantities which are relative to a fixed coordinate system are distinguished by using the subscript,  $B$  on those relative to the moving coordinates system.

$$\mathbf{a} = \mathbf{v}' = \mathbf{R}'' + x''\mathbf{i} + y''\mathbf{j} + z''\mathbf{k} + \overbrace{x\mathbf{i}' + y\mathbf{j}' + z\mathbf{k}'}^{\boldsymbol{\Omega} \times \mathbf{v}_B} + \boldsymbol{\Omega}' \times \mathbf{r}_B$$

$$\begin{aligned}
& + \Omega \times \left( \overbrace{x'\mathbf{i} + y'\mathbf{j} + z'\mathbf{k}}^{\mathbf{v}_B} + \overbrace{xi' + yj' + zk'}^{\Omega \times \mathbf{r}_B(t)} \right) \\
& = \mathbf{R}'' + \mathbf{a}_B + \Omega' \times \mathbf{r}_B + 2\Omega \times \mathbf{v}_B + \Omega \times (\Omega \times \mathbf{r}_B).
\end{aligned}$$

The acceleration  $\mathbf{a}_B$  is that perceived by an observer who is moving with the moving coordinate system and for whom the moving coordinate system is fixed. The term  $\Omega \times (\Omega \times \mathbf{r}_B)$  is called the centripetal acceleration. Solving for  $\mathbf{a}_B$ ,

$$\mathbf{a}_B = \mathbf{a} - \mathbf{R}'' - \Omega' \times \mathbf{r}_B - 2\Omega \times \mathbf{v}_B - \Omega \times (\Omega \times \mathbf{r}_B). \quad (19.30)$$

Here the term  $-(\Omega \times (\Omega \times \mathbf{r}_B))$  is called the centrifugal acceleration, it being an acceleration felt by the observer relative to the moving coordinate system which he regards as fixed, and the term  $-2\Omega \times \mathbf{v}_B$  is called the Coriolis acceleration, an acceleration experienced by the observer as he moves relative to the moving coordinate system. The mass multiplied by the Coriolis acceleration defines the Coriolis force.

There is a ride found in some amusement parks in which the victims stand next to a circular wall covered with a carpet or some rough material. Then the whole circular room begins to revolve faster and faster. At some point, the bottom drops out and the victims are held in place by friction. The force they feel is called centrifugal force and it causes centrifugal acceleration. It is not necessary to move relative to coordinates fixed with the revolving wall in order to feel this force and it is pretty predictable. However, if the nauseated victim moves relative to the rotating wall, he will feel the effects of the Coriolis force and this force is really strange. The difference between these forces is that the Coriolis force is caused by movement relative to the moving coordinate system and the centrifugal force is not.

### 19.5.2 The Coriolis Acceleration On The Rotating Earth

Now consider the earth. Let  $\mathbf{i}^*, \mathbf{j}^*, \mathbf{k}^*$ , be the usual basis vectors attached to the rotating earth. Thus  $\mathbf{k}^*$  is fixed in space with  $\mathbf{k}^*$  pointing in the direction of the north pole from the center of the earth while  $\mathbf{i}^*$  and  $\mathbf{j}^*$  point to fixed points on the surface of the earth. Thus  $\mathbf{i}^*$  and  $\mathbf{j}^*$  depend on  $t$  while  $\mathbf{k}^*$  does not. Let  $\mathbf{i}, \mathbf{j}, \mathbf{k}$  be the unit vectors described earlier with  $\mathbf{i}$  pointing South,  $\mathbf{j}$  pointing East, and  $\mathbf{k}$  pointing away from the center of the earth at some point of the rotating earth's surface,  $\mathbf{p}$ . Letting  $\mathbf{R}(t)$  be the position vector of the point  $\mathbf{p}$ , from the center of the earth, observe the coordinates of  $\mathbf{R}(t)$  are constant with respect to  $\mathbf{i}(t), \mathbf{j}(t), \mathbf{k}(t)$ . Also, since the earth rotates from West to East and the speed of a point on the surface of the earth relative to an observer fixed in space is  $\omega |\mathbf{R}| \sin \phi$  where  $\omega$  is the angular speed of the earth about an axis through the poles, it follows from the geometric definition of the cross product that

$$\mathbf{R}' = \omega \mathbf{k}^* \times \mathbf{R}$$

Therefore,  $\Omega = \omega \mathbf{k}^*$  and so

$$\mathbf{R}'' = \overbrace{\Omega' \times \mathbf{R}}^{=0} + \Omega \times \mathbf{R}' = \Omega \times (\Omega \times \mathbf{R})$$

since  $\Omega$  does not depend on  $t$ . Formula (19.30) implies

$$\mathbf{a}_B = \mathbf{a} - \Omega \times (\Omega \times \mathbf{R}) - 2\Omega \times \mathbf{v}_B - \Omega \times (\Omega \times \mathbf{r}_B). \quad (19.31)$$



In this formula, you can totally ignore the term  $\boldsymbol{\Omega} \times (\boldsymbol{\Omega} \times \mathbf{r}_B)$  because it is so small whenever you are considering motion near some point on the earth's surface. To see this, note

$\omega \overbrace{(24)(3600)}^{\text{seconds in a day}} = 2\pi$ , and so  $\omega = 7.2722 \times 10^{-5}$  in radians per second. If you are using seconds to measure time and feet to measure distance, this term is therefore, no larger than

$$(7.2722 \times 10^{-5})^2 |\mathbf{r}_B|.$$

Clearly this is not worth considering in the presence of the acceleration due to gravity which is approximately 32 feet per second squared near the surface of the earth.

If the acceleration  $\mathbf{a}$ , is due to gravity, then

$$\begin{aligned} \mathbf{a}_B &= \mathbf{a} - \boldsymbol{\Omega} \times (\boldsymbol{\Omega} \times \mathbf{R}) - 2\boldsymbol{\Omega} \times \mathbf{v}_B = \\ &= \overbrace{-\frac{GM(\mathbf{R} + \mathbf{r}_B)}{|\mathbf{R} + \mathbf{r}_B|^3} - \boldsymbol{\Omega} \times (\boldsymbol{\Omega} \times \mathbf{R}) - 2\boldsymbol{\Omega} \times \mathbf{v}_B}^{\equiv \mathbf{g}} \equiv \mathbf{g} - 2\boldsymbol{\Omega} \times \mathbf{v}_B. \end{aligned}$$

Note that

$$\boldsymbol{\Omega} \times (\boldsymbol{\Omega} \times \mathbf{R}) = (\boldsymbol{\Omega} \cdot \mathbf{R}) \boldsymbol{\Omega} - |\boldsymbol{\Omega}|^2 \mathbf{R}$$

and so  $\mathbf{g}$ , the acceleration relative to the moving coordinate system on the earth is not directed exactly toward the center of the earth except at the poles and at the equator, although the components of acceleration which are in other directions are very small when compared with the acceleration due to the force of gravity and are often neglected. Therefore, if the only force acting on an object is due to gravity, the following formula describes the acceleration relative to a coordinate system moving with the earth's surface.

$$\mathbf{a}_B = \mathbf{g} - 2(\boldsymbol{\Omega} \times \mathbf{v}_B)$$

While the vector,  $\boldsymbol{\Omega}$  is quite small, if the relative velocity,  $\mathbf{v}_B$  is large, the Coriolis acceleration could be significant. This is described in terms of the vectors  $\mathbf{i}(t), \mathbf{j}(t), \mathbf{k}(t)$  next.

Letting  $(\rho, \theta, \phi)$  be the usual spherical coordinates of the point  $\mathbf{p}(t)$  on the surface taken with respect to  $\mathbf{i}^*, \mathbf{j}^*, \mathbf{k}^*$  the usual way with  $\phi$  the polar angle, it follows the  $\mathbf{i}^*, \mathbf{j}^*, \mathbf{k}^*$  coordinates of this point are

$$\begin{pmatrix} \rho \sin(\phi) \cos(\theta) \\ \rho \sin(\phi) \sin(\theta) \\ \rho \cos(\phi) \end{pmatrix}.$$

It follows,

$$\mathbf{i} = \cos(\phi) \cos(\theta) \mathbf{i}^* + \cos(\phi) \sin(\theta) \mathbf{j}^* - \sin(\phi) \mathbf{k}^*$$

$$\mathbf{j} = -\sin(\theta) \mathbf{i}^* + \cos(\theta) \mathbf{j}^* + 0 \mathbf{k}^*$$

and

$$\mathbf{k} = \sin(\phi) \cos(\theta) \mathbf{i}^* + \sin(\phi) \sin(\theta) \mathbf{j}^* + \cos(\phi) \mathbf{k}^*.$$

It is necessary to obtain  $\mathbf{k}^*$  in terms of the vectors,  $\mathbf{i}, \mathbf{j}, \mathbf{k}$ . Thus the following equation needs to be solved for  $a, b, c$  to find  $\mathbf{k}^* = a\mathbf{i} + b\mathbf{j} + c\mathbf{k}$

$$\overbrace{\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}}^{\mathbf{k}^*} = \begin{pmatrix} \cos(\phi) \cos(\theta) & -\sin(\theta) & \sin(\phi) \cos(\theta) \\ \cos(\phi) \sin(\theta) & \cos(\theta) & \sin(\phi) \sin(\theta) \\ -\sin(\phi) & 0 & \cos(\phi) \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} \quad (19.32)$$

The first column is  $\mathbf{i}$ , the second is  $\mathbf{j}$  and the third is  $\mathbf{k}$  in the above matrix. The solution is  $a = -\sin(\phi)$ ,  $b = 0$ , and  $c = \cos(\phi)$ .

Now the Coriolis acceleration on the earth equals

$$2(\boldsymbol{\Omega} \times \mathbf{v}_B) = 2\omega \left( \overbrace{-\sin(\phi)\mathbf{i} + 0\mathbf{j} + \cos(\phi)\mathbf{k}}^{\mathbf{k}^*} \right) \times (x'\mathbf{i} + y'\mathbf{j} + z'\mathbf{k}).$$

This equals

$$2\omega [(-y' \cos \phi)\mathbf{i} + (x' \cos \phi + z' \sin \phi)\mathbf{j} - (y' \sin \phi)\mathbf{k}]. \quad (19.33)$$

Remember  $\phi$  is fixed and pertains to the fixed point,  $\mathbf{p}(t)$  on the earth's surface. Therefore, if the acceleration,  $\mathbf{a}$  is due to gravity,

$$\mathbf{a}_B = \mathbf{g} - 2\omega [(-y' \cos \phi)\mathbf{i} + (x' \cos \phi + z' \sin \phi)\mathbf{j} - (y' \sin \phi)\mathbf{k}]$$

where  $\mathbf{g} = -\frac{GM(\mathbf{R} + \mathbf{r}_B)}{|\mathbf{R} + \mathbf{r}_B|^3} - \boldsymbol{\Omega} \times (\boldsymbol{\Omega} \times \mathbf{R})$  as explained above. The term  $\boldsymbol{\Omega} \times (\boldsymbol{\Omega} \times \mathbf{R})$  is pretty small and so it will be neglected. However, the Coriolis force will not be neglected.

**Example 19.5.5** Suppose a rock is dropped from a tall building. Where will it strike?

Assume  $\mathbf{a} = -g\mathbf{k}$  and the  $\mathbf{j}$  component of  $\mathbf{a}_B$  is approximately

$$-2\omega (x' \cos \phi + z' \sin \phi).$$

The dominant term in this expression is clearly the second one because  $x'$  will be small. Also, the  $\mathbf{i}$  and  $\mathbf{k}$  contributions will be very small. Therefore, the following equation is descriptive of the situation.

$$\mathbf{a}_B = -g\mathbf{k} - 2z'\omega \sin \phi \mathbf{j}.$$

$z' = -gt$  approximately. Therefore, considering the  $\mathbf{j}$  component, this is

$$2gt\omega \sin \phi.$$

Two integrations give  $(\omega g t^3 / 3) \sin \phi$  for the  $\mathbf{j}$  component of the relative displacement at time  $t$ .

This shows the rock does not fall directly towards the center of the earth as expected but slightly to the east.

**Example 19.5.6** In 1851 Foucault set a pendulum vibrating and observed the earth rotate out from under it. It was a very long pendulum with a heavy weight at the end so that it would vibrate for a long time without stopping<sup>1</sup>. This is what allowed him to observe the earth rotate out from under it. Clearly such a pendulum will take 24 hours for the plane of vibration to appear to make one complete revolution at the north pole. It is also reasonable to expect that no such observed rotation would take place on the equator. Is it possible to predict what will take place at various latitudes?

Using (19.33), in (19.31),

$$\mathbf{a}_B = \mathbf{a} - \boldsymbol{\Omega} \times (\boldsymbol{\Omega} \times \mathbf{R}) - 2\omega [(-y' \cos \phi)\mathbf{i} + (x' \cos \phi + z' \sin \phi)\mathbf{j} - (y' \sin \phi)\mathbf{k}].$$

<sup>1</sup>There is such a pendulum in the Eyring building at BYU and to keep people from touching it, there is a little sign which says Warning! 1000 ohms.

Neglecting the small term,  $\boldsymbol{\Omega} \times (\boldsymbol{\Omega} \times \mathbf{R})$ , this becomes

$$= -g\mathbf{k} + \mathbf{T}/m - 2\omega [(-y' \cos \phi) \mathbf{i} + (x' \cos \phi + z' \sin \phi) \mathbf{j} - (y' \sin \phi) \mathbf{k}]$$

where  $\mathbf{T}$ , the tension in the string of the pendulum, is directed towards the point at which the pendulum is supported, and  $m$  is the mass of the pendulum bob. The pendulum can be thought of as the position vector from  $(0, 0, l)$  to the surface of the sphere  $x^2 + y^2 + (z - l)^2 = l^2$ . Therefore,

$$\mathbf{T} = -T \frac{x}{l} \mathbf{i} - T \frac{y}{l} \mathbf{j} + T \frac{l - z}{l} \mathbf{k}$$

and consequently, the differential equations of relative motion are

$$x'' = -T \frac{x}{ml} + 2\omega y' \cos \phi$$

$$y'' = -T \frac{y}{ml} - 2\omega (x' \cos \phi + z' \sin \phi)$$

and

$$z'' = T \frac{l - z}{ml} - g + 2\omega y' \sin \phi.$$

If the vibrations of the pendulum are small so that for practical purposes,  $z'' = z = 0$ , the last equation may be solved for  $T$  to get

$$gm - 2\omega y' \sin(\phi) m = T.$$

Therefore, the first two equations become

$$x'' = -(gm - 2\omega m y' \sin \phi) \frac{x}{ml} + 2\omega y' \cos \phi$$

and

$$y'' = -(gm - 2\omega m y' \sin \phi) \frac{y}{ml} - 2\omega (x' \cos \phi + z' \sin \phi).$$

All terms of the form  $xy'$  or  $y'y$  can be neglected because it is assumed  $x$  and  $y$  remain small. Also, the pendulum is assumed to be long with a heavy weight so that  $x'$  and  $y'$  are also small. With these simplifying assumptions, the equations of motion become

$$x'' + g \frac{x}{l} = 2\omega y' \cos \phi$$

and

$$y'' + g \frac{y}{l} = -2\omega x' \cos \phi.$$

These equations are of the form

$$x'' + a^2 x = by', \quad y'' + a^2 y = -bx' \quad (19.34)$$

where  $a^2 = \frac{g}{l}$  and  $b = 2\omega \cos \phi$ . Then it is fairly tedious but routine to verify that for each constant,  $c$ ,

$$x = c \sin\left(\frac{bt}{2}\right) \sin\left(\frac{\sqrt{b^2 + 4a^2}}{2} t\right), \quad y = c \cos\left(\frac{bt}{2}\right) \sin\left(\frac{\sqrt{b^2 + 4a^2}}{2} t\right) \quad (19.35)$$

yields a solution to (19.34) along with the initial conditions,

$$x(0) = 0, y(0) = 0, x'(0) = 0, y'(0) = \frac{c\sqrt{b^2 + 4a^2}}{2}. \quad (19.36)$$

It is clear from experiments with the pendulum that the earth does indeed rotate out from under it causing the plane of vibration of the pendulum to appear to rotate. The purpose of this discussion is not to establish these self evident facts but to predict how long it takes for the plane of vibration to make one revolution. Therefore, there will be some instant in time at which the pendulum will be vibrating in a plane determined by  $\mathbf{k}$  and  $\mathbf{j}$ . (Recall  $\mathbf{k}$  points away from the center of the earth and  $\mathbf{j}$  points East. ) At this instant in time, defined as  $t = 0$ , the conditions of (19.36) will hold for some value of  $c$  and so the solution to (19.34) having these initial conditions will be those of (19.35) by uniqueness of the initial value problem. Writing these solutions differently,

$$\begin{pmatrix} x(t) \\ y(t) \end{pmatrix} = c \begin{pmatrix} \sin\left(\frac{bt}{2}\right) \\ \cos\left(\frac{bt}{2}\right) \end{pmatrix} \sin\left(\frac{\sqrt{b^2 + 4a^2}}{2}t\right)$$

This is very interesting! The vector,  $c \begin{pmatrix} \sin\left(\frac{bt}{2}\right) \\ \cos\left(\frac{bt}{2}\right) \end{pmatrix}$  always has magnitude equal to  $|c|$  but its direction changes very slowly because  $b$  is very small. The plane of vibration is determined by this vector and the vector  $\mathbf{k}$ . The term  $\sin\left(\frac{\sqrt{b^2 + 4a^2}}{2}t\right)$  changes relatively fast and takes values between  $-1$  and  $1$ . This is what describes the actual observed vibrations of the pendulum. Thus the plane of vibration will have made one complete revolution when  $t = T$  for

$$\frac{bT}{2} \equiv 2\pi.$$

Therefore, the time it takes for the earth to turn out from under the pendulum is

$$T = \frac{4\pi}{2\omega \cos \phi} = \frac{2\pi}{\omega} \sec \phi.$$

Since  $\omega$  is the angular speed of the rotating earth, it follows  $\omega = \frac{2\pi}{24} = \frac{\pi}{12}$  in radians per hour. Therefore, the above formula implies

$$T = 24 \sec \phi.$$

I think this is really amazing. You could actually determine latitude, not by taking readings with instruments using the North Star but by doing an experiment with a big pendulum. You would set it vibrating, observe  $T$  in hours, and then solve the above equation for  $\phi$ . Also note the pendulum would not appear to change its plane of vibration at the equator because  $\lim_{\phi \rightarrow \pi/2} \sec \phi = \infty$ .

The Coriolis acceleration is also responsible for the phenomenon of the next example.

**Example 19.5.7** *It is known that low pressure areas rotate counterclockwise as seen from above in the Northern hemisphere but clockwise in the Southern hemisphere. Why?*

Neglect accelerations other than the Coriolis acceleration and the following acceleration which comes from an assumption that the point  $\mathbf{p}(t)$  is the location of the lowest pressure.

$$\mathbf{a} = -a(r_B) \mathbf{r}_B$$

where  $r_B = r$  will denote the distance from the fixed point  $\mathbf{p}(t)$  on the earth's surface which is also the lowest pressure point. Of course the situation could be more complicated but this will suffice to explain the above question. Then the acceleration observed by a person on the earth relative to the apparently fixed vectors,  $\mathbf{i}, \mathbf{k}, \mathbf{j}$ , is

$$\mathbf{a}_B = -a(r_B)(x\mathbf{i} + y\mathbf{j} + z\mathbf{k}) - 2\omega[-y'\cos(\phi)\mathbf{i} + (x'\cos(\phi) + z'\sin(\phi))\mathbf{j} - (y'\sin(\phi)\mathbf{k})]$$

Therefore, one obtains some differential equations from  $\mathbf{a}_B = x''\mathbf{i} + y''\mathbf{j} + z''\mathbf{k}$  by matching the components. These are

$$\begin{aligned}x'' + a(r_B)x &= 2\omega y' \cos \phi \\y'' + a(r_B)y &= -2\omega x' \cos \phi - 2\omega z' \sin(\phi) \\z'' + a(r_B)z &= 2\omega y' \sin \phi\end{aligned}$$

Now remember, the vectors,  $\mathbf{i}, \mathbf{j}, \mathbf{k}$  are fixed relative to the earth and so are constant vectors. Therefore, from the properties of the determinant and the above differential equations,

$$\begin{aligned}(\mathbf{r}'_B \times \mathbf{r}_B)' &= \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ x' & y' & z' \\ x & y & z \end{vmatrix}' = \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ x'' & y'' & z'' \\ x & y & z \end{vmatrix} \\ &= \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ -a(r_B)x + 2\omega y' \cos \phi & -a(r_B)y - 2\omega x' \cos \phi - 2\omega z' \sin(\phi) & -a(r_B)z + 2\omega y' \sin \phi \\ x & y & z \end{vmatrix}\end{aligned}$$

Then the  $\mathbf{k}^{th}$  component of this cross product equals

$$\omega \cos(\phi) (y^2 + x^2)' + 2\omega x z' \sin(\phi).$$

The first term will be negative because it is assumed  $\mathbf{p}(t)$  is the location of low pressure causing  $y^2 + x^2$  to be a decreasing function. If it is assumed there is not a substantial motion in the  $\mathbf{k}$  direction, so that  $z$  is fairly constant and the last term can be neglected, then the  $\mathbf{k}^{th}$  component of  $(\mathbf{r}'_B \times \mathbf{r}_B)'$  is negative provided  $\phi \in (0, \frac{\pi}{2})$  and positive if  $\phi \in (\frac{\pi}{2}, \pi)$ . Beginning with a point at rest, this implies  $\mathbf{r}'_B \times \mathbf{r}_B = \mathbf{0}$  initially and then the above implies its  $\mathbf{k}^{th}$  component is negative in the upper hemisphere when  $\phi < \pi/2$  and positive in the lower hemisphere when  $\phi > \pi/2$ . Using the right hand and the geometric definition of the cross product, this shows clockwise rotation in the lower hemisphere and counter clockwise rotation in the upper hemisphere.

Note also that as  $\phi$  gets close to  $\pi/2$  near the equator, the above reasoning tends to break down because  $\cos(\phi)$  becomes close to zero. Therefore, the motion towards the low pressure has to be more pronounced in comparison with the motion in the  $\mathbf{k}$  direction in order to draw this conclusion.

## 19.6 Exercises

1. Remember the Coriolis force was  $2\boldsymbol{\Omega} \times \mathbf{v}_B$  where  $\boldsymbol{\Omega}$  was a particular vector which came from the matrix,  $Q(t)$  as described above. Show that

$$Q(t) = \begin{pmatrix} \mathbf{i}(t) \cdot \mathbf{i}(t_0) & \mathbf{j}(t) \cdot \mathbf{i}(t_0) & \mathbf{k}(t) \cdot \mathbf{i}(t_0) \\ \mathbf{i}(t) \cdot \mathbf{j}(t_0) & \mathbf{j}(t) \cdot \mathbf{j}(t_0) & \mathbf{k}(t) \cdot \mathbf{j}(t_0) \\ \mathbf{i}(t) \cdot \mathbf{k}(t_0) & \mathbf{j}(t) \cdot \mathbf{k}(t_0) & \mathbf{k}(t) \cdot \mathbf{k}(t_0) \end{pmatrix}.$$

There will be no Coriolis force exactly when  $\boldsymbol{\Omega} = \mathbf{0}$  which corresponds to  $Q'(t) = 0$ . When will  $Q'(t) = 0$ ?

2. An illustration used in many beginning physics books is that of firing a rifle horizontally and dropping an identical bullet from the same height above the perfectly flat ground followed by an assertion that the two bullets will hit the ground at exactly the same time. Is this true on the rotating earth assuming the experiment

takes place over a large perfectly flat field so the curvature of the earth is not an issue? Explain. What other irregularities will occur? Recall the Coriolis force is  $2\omega [(-y' \cos \phi) \mathbf{i} + (x' \cos \phi + z' \sin \phi) \mathbf{j} - (y' \sin \phi) \mathbf{k}]$  where  $\mathbf{k}$  points away from the center of the earth,  $\mathbf{j}$  points East, and  $\mathbf{i}$  points South.

## 19.7 Determinants

Let  $A$  be an  $n \times n$  matrix. The determinant of  $A$ , denoted as  $\det(A)$  is a number. If the matrix is a  $2 \times 2$  matrix, this number is very easy to find.

**Definition 19.7.1** Let  $A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$ . Then

$$\det(A) \equiv ad - cb.$$

The determinant is also often denoted by enclosing the matrix with two vertical lines. Thus

$$\det \begin{pmatrix} a & b \\ c & d \end{pmatrix} = \begin{vmatrix} a & b \\ c & d \end{vmatrix}.$$

**Example 19.7.2** Find  $\det \begin{pmatrix} 2 & 4 \\ -1 & 6 \end{pmatrix}$ .

From the definition this is just  $(2)(6) - (-1)(4) = 16$ .

Having defined what is meant by the determinant of a  $2 \times 2$  matrix, what about a  $3 \times 3$  matrix?

**Example 19.7.3** Find the determinant of

$$\begin{pmatrix} 1 & 2 & 3 \\ 4 & 3 & 2 \\ 3 & 2 & 1 \end{pmatrix}.$$

Here is how it is done by “expanding along the first column”.

$$(-1)^{1+1} 1 \begin{vmatrix} 3 & 2 \\ 2 & 1 \end{vmatrix} + (-1)^{2+1} 4 \begin{vmatrix} 2 & 3 \\ 2 & 1 \end{vmatrix} + (-1)^{3+1} 3 \begin{vmatrix} 2 & 3 \\ 3 & 2 \end{vmatrix} = 0.$$

What is going on here? Take the 1 in the upper left corner and cross out the row and the column containing the 1. Then take the determinant of the resulting  $2 \times 2$  matrix. Now multiply this determinant by 1 and then multiply by  $(-1)^{1+1}$  because this 1 is in the first row and the first column. This gives the first term in the above sum. Now go to the 4. Cross out the row and the column which contain 4 and take the determinant of the  $2 \times 2$  matrix which remains. Multiply this by 4 and then by  $(-1)^{2+1}$  because the 4 is in the first column and the second row. Finally consider the 3 on the bottom of the first column. Cross out the row and column containing this 3 and take the determinant of what is left. Then multiply this by 3 and by  $(-1)^{3+1}$  because this 3 is in the third row and the first column. This is the pattern used to evaluate the determinant by expansion along the first column.

You could also expand the determinant along the second row as follows.

$$(-1)^{2+1} 4 \begin{vmatrix} 2 & 3 \\ 2 & 1 \end{vmatrix} + (-1)^{2+2} 3 \begin{vmatrix} 1 & 3 \\ 3 & 1 \end{vmatrix} + (-1)^{2+3} 2 \begin{vmatrix} 1 & 2 \\ 3 & 2 \end{vmatrix} = 0.$$

It follows exactly the same pattern and you see it gave the same answer. You pick a row or column and corresponding to each number in that row or column, you cross out the row

and column containing it, take the determinant of what is left, multiply this by the number and by  $(-1)^{i+j}$  assuming the number is in the  $i^{\text{th}}$  row and the  $j^{\text{th}}$  column. Then adding these gives the value of the determinant.

What about a  $4 \times 4$  matrix?

**Example 19.7.4** Find  $\det(A)$  where

$$A = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 5 & 4 & 2 & 3 \\ 1 & 3 & 4 & 5 \\ 3 & 4 & 3 & 2 \end{pmatrix}$$

As in the case of a  $3 \times 3$  matrix, you can expand this along any row or column. Lets pick the third column.  $\det(A) =$

$$3(-1)^{1+3} \begin{vmatrix} 5 & 4 & 3 \\ 1 & 3 & 5 \\ 3 & 4 & 2 \end{vmatrix} + 2(-1)^{2+3} \begin{vmatrix} 1 & 2 & 4 \\ 1 & 3 & 5 \\ 3 & 4 & 2 \end{vmatrix} +$$

$$4(-1)^{3+3} \begin{vmatrix} 1 & 2 & 4 \\ 5 & 4 & 3 \\ 3 & 4 & 2 \end{vmatrix} + 3(-1)^{4+3} \begin{vmatrix} 1 & 2 & 4 \\ 5 & 4 & 3 \\ 1 & 3 & 5 \end{vmatrix}.$$

Now you know how to expand each of these  $3 \times 3$  matrices along a row or a column. If you do so, you will get  $-12$  assuming you make no mistakes. You could expand this matrix along any row or any column and assuming you make no mistakes, you will always get the same thing which is defined to be the determinant of the matrix,  $A$ . This method of evaluating a determinant by expanding along a row or a column is called the method of Laplace expansion.

Note that each of the four terms above involves three terms consisting of determinants of  $2 \times 2$  matrices and each of these will need 2 terms. Therefore, there will be  $4 \times 3 \times 2 = 24$  terms to evaluate in order to find the determinant using the method of Laplace expansion. Suppose now you have a  $10 \times 10$  matrix. I hope you see that from the above pattern there will be  $10! = 3,628,800$  terms involved in the evaluation of such a determinant by Laplace expansion along a row or column. This is a lot of terms.

In addition to the difficulties just discussed, I think you should regard the above claim that you always get the same answer by picking any row or column with considerable skepticism. It is incredible and not at all obvious. However, it requires a little effort to establish it. This is done in the section on the theory of the determinant which follows. The above examples motivate the following incredible theorem and definition.

**Definition 19.7.5** Let  $A = (a_{ij})$  be an  $n \times n$  matrix. Then a new matrix called the cofactor matrix,  $\text{cof}(A)$  is defined by  $\text{cof}(A) = (c_{ij})$  where to obtain  $c_{ij}$  delete the  $i^{\text{th}}$  row and the  $j^{\text{th}}$  column of  $A$ , take the determinant of the  $(n-1) \times (n-1)$  matrix which results, (This is called the  $ij^{\text{th}}$  minor of  $A$ .) and then multiply this number by  $(-1)^{i+j}$ . To make the formulas easier to remember,  $\text{cof}(A)_{ij}$  will denote the  $ij^{\text{th}}$  entry of the cofactor matrix.

**Theorem 19.7.6** Let  $A$  be an  $n \times n$  matrix where  $n \geq 2$ . Then

$$\det(A) = \sum_{j=1}^n a_{ij} \text{cof}(A)_{ij} = \sum_{i=1}^n a_{ij} \text{cof}(A)_{ij}. \quad (19.37)$$

The first formula consists of expanding the determinant along the  $i^{\text{th}}$  row and the second expands the determinant along the  $j^{\text{th}}$  column.

Notwithstanding the difficulties involved in using the method of Laplace expansion, certain types of matrices are very easy to deal with.

**Definition 19.7.7** A matrix  $M$ , is upper triangular if  $M_{ij} = 0$  whenever  $i > j$ . Thus such a matrix equals zero below the main diagonal, the entries of the form  $M_{ii}$ , as shown.

$$\begin{pmatrix} * & * & \cdots & * \\ 0 & * & \ddots & \vdots \\ \vdots & \ddots & \ddots & * \\ 0 & \cdots & 0 & * \end{pmatrix}$$

A lower triangular matrix is defined similarly as a matrix for which all entries above the main diagonal are equal to zero.

You should verify the following using the above theorem on Laplace expansion.

**Corollary 19.7.8** Let  $M$  be an upper (lower) triangular matrix. Then  $\det(M)$  is obtained by taking the product of the entries on the main diagonal.

**Example 19.7.9** Let

$$A = \begin{pmatrix} 1 & 2 & 3 & 77 \\ 0 & 2 & 6 & 7 \\ 0 & 0 & 3 & 33.7 \\ 0 & 0 & 0 & -1 \end{pmatrix}$$

Find  $\det(A)$ .

From the above corollary, it suffices to take the product of the diagonal elements. Thus  $\det(A) = 1 \times 2 \times 3 \times -1 = -6$ . Without using the corollary, you could expand along the first column. This gives

$$1 \begin{vmatrix} 2 & 6 & 7 \\ 0 & 3 & 33.7 \\ 0 & 0 & -1 \end{vmatrix}$$

and now expand this along the first column to get this equals

$$1 \times 2 \times \begin{vmatrix} 3 & 33.7 \\ 0 & -1 \end{vmatrix}$$

Next expand the last along the first column which reduces to the product of the main diagonal elements as claimed. This example also demonstrates why the above corollary is true.

There are many properties satisfied by determinants. Some of the most important are listed in the following theorem.

**Theorem 19.7.10** If two rows or two columns in an  $n \times n$  matrix,  $A$ , are switched, the determinant of the resulting matrix equals  $(-1)$  times the determinant of the original matrix. If  $A$  is an  $n \times n$  matrix in which two rows are equal or two columns are equal then  $\det(A) = 0$ . Suppose the  $i^{\text{th}}$  row of  $A$  equals  $(xa_1 + yb_1, \dots, xa_n + yb_n)$ . Then

$$\det(A) = x \det(A_1) + y \det(A_2)$$

where the  $i^{\text{th}}$  row of  $A_1$  is  $(a_1, \dots, a_n)$  and the  $i^{\text{th}}$  row of  $A_2$  is  $(b_1, \dots, b_n)$ , all other rows of  $A_1$  and  $A_2$  coinciding with those of  $A$ . In other words,  $\det$  is a linear function of each



row  $A$ . The same is true with the word “row” replaced with the word “column”. In addition to this, if  $A$  and  $B$  are  $n \times n$  matrices, then

$$\det(AB) = \det(A) \det(B),$$

and if  $A$  is an  $n \times n$  matrix, then

$$\det(A) = \det(A^T).$$

This theorem implies the following corollary which gives a way to find determinants. As I pointed out above, the method of Laplace expansion will not be practical for any matrix of large size.

**Corollary 19.7.11** *Let  $A$  be an  $n \times n$  matrix and let  $B$  be the matrix obtained by replacing the  $i^{\text{th}}$  row (column) of  $A$  with the sum of the  $i^{\text{th}}$  row (column) added to a multiple of another row (column). Then  $\det(A) = \det(B)$ . If  $B$  is the matrix obtained from  $A$  by replacing the  $i^{\text{th}}$  row (column) of  $A$  by a times the  $i^{\text{th}}$  row (column) then  $a \det(A) = \det(B)$ .*

Here is an example which shows how to use this corollary to find a determinant.

**Example 19.7.12** *Find the determinant of the matrix,*

$$A = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 5 & 1 & 2 & 3 \\ 4 & 5 & 4 & 3 \\ 2 & 2 & -4 & 5 \end{pmatrix}$$

Replace the second row by  $(-5)$  times the first row added to it. Then replace the third row by  $(-4)$  times the first row added to it. Finally, replace the fourth row by  $(-2)$  times the first row added to it. This yields the matrix,

$$B = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 0 & -9 & -13 & -17 \\ 0 & -3 & -8 & -13 \\ 0 & -2 & -10 & -3 \end{pmatrix}$$

and from the above corollary, it has the same determinant as  $A$ . Now using the corollary some more,  $\det(B) = \left(\frac{-1}{3}\right) \det(C)$  where

$$C = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 0 & 0 & 11 & 22 \\ 0 & -3 & -8 & -13 \\ 0 & 6 & 30 & 9 \end{pmatrix}.$$

The second row was replaced by  $(-3)$  times the third row added to the second row and then the last row was multiplied by  $(-3)$ . Now replace the last row with 2 times the third added to it and then switch the third and second rows. Then  $\det(C) = -\det(D)$  where

$$D = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 0 & -3 & -8 & -13 \\ 0 & 0 & 11 & 22 \\ 0 & 0 & 14 & -17 \end{pmatrix}$$

You could do more row operations or you could note that this can be easily expanded along the first column followed by expanding the  $3 \times 3$  matrix which results along its first column. Thus

$$\det(D) = 1(-3) \begin{vmatrix} 11 & 22 \\ 14 & -17 \end{vmatrix} = 1485$$

and so  $\det(C) = -1485$  and  $\det(A) = \det(B) = \left(\frac{-1}{3}\right)(-1485) = 495$ .

The theorem about expanding a matrix along any row or column also provides a way to give a formula for the inverse of a matrix. Recall the definition of the inverse of a matrix in Definition 19.1.13 on Page 417.

**Theorem 19.7.13**  $A^{-1}$  exists if and only if  $\det(A) \neq 0$ . If  $\det(A) \neq 0$ , then  $A^{-1} = (a_{ij}^{-1})$  where

$$a_{ij}^{-1} = \det(A)^{-1} \operatorname{cof}(A)_{ji}$$

for  $\operatorname{cof}(A)_{ij}$  the  $ij^{\text{th}}$  cofactor of  $A$ .

**Proof:** By Theorem 19.7.6 and letting  $(a_{ir}) = A$ , if  $\det(A) \neq 0$ ,

$$\sum_{i=1}^n a_{ir} \operatorname{cof}(A)_{ir} \det(A)^{-1} = \det(A) \det(A)^{-1} = 1.$$

Now consider

$$\sum_{i=1}^n a_{ir} \operatorname{cof}(A)_{ik} \det(A)^{-1}$$

when  $k \neq r$ . Replace the  $k^{\text{th}}$  column with the  $r^{\text{th}}$  column to obtain a matrix,  $B_k$  whose determinant equals zero by Theorem 19.7.10. However, expanding this matrix along the  $k^{\text{th}}$  column yields

$$0 = \det(B_k) \det(A)^{-1} = \sum_{i=1}^n a_{ir} \operatorname{cof}(A)_{ik} \det(A)^{-1}$$

Summarizing,

$$\sum_{i=1}^n a_{ir} \operatorname{cof}(A)_{ik} \det(A)^{-1} = \delta_{rk}.$$

Now

$$\sum_{i=1}^n a_{ir} \operatorname{cof}(A)_{ik} = \sum_{i=1}^n a_{ir} \operatorname{cof}(A)_{ki}^T$$

which is the  $kr^{\text{th}}$  entry of  $\operatorname{cof}(A)^T A$ . Therefore,

$$\frac{\operatorname{cof}(A)^T}{\det(A)} A = I. \quad (19.38)$$

Using the other formula in Theorem 19.7.6, and similar reasoning,

$$\sum_{j=1}^n a_{rj} \operatorname{cof}(A)_{kj} \det(A)^{-1} = \delta_{rk}$$

Now

$$\sum_{j=1}^n a_{rj} \operatorname{cof}(A)_{kj} = \sum_{j=1}^n a_{rj} \operatorname{cof}(A)_{jk}^T$$

which is the  $rk^{th}$  entry of  $A \operatorname{cof}(A)^T$ . Therefore,

$$A \frac{\operatorname{cof}(A)^T}{\det(A)} = I, \quad (19.39)$$

and it follows from (19.38) and (19.39) that  $A^{-1} = (a_{ij}^{-1})$ , where

$$a_{ij}^{-1} = \operatorname{cof}(A)_{ji} \det(A)^{-1}.$$

In other words,

$$A^{-1} = \frac{\operatorname{cof}(A)^T}{\det(A)}.$$

Now suppose  $A^{-1}$  exists. Then by Theorem 19.7.10,

$$1 = \det(I) = \det(AA^{-1}) = \det(A) \det(A^{-1})$$

so  $\det(A) \neq 0$ . This proves the theorem.

Theorem 19.7.13 says that to find the inverse, take the transpose of the cofactor matrix and divide by the determinant. The transpose of the cofactor matrix is called the adjugate or sometimes the classical adjoint of the matrix  $A$ . It is an abomination to call it the adjoint although you do sometimes see it referred to in this way. In words,  $A^{-1}$  is equal to one over the determinant of  $A$  times the adjugate matrix of  $A$ .

**Example 19.7.14** Find the inverse of the matrix,

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 3 & 0 & 1 \\ 1 & 2 & 1 \end{pmatrix}$$

First find the determinant of this matrix. Using Corollary 19.7.11 on Page 441, the determinant of this matrix equals the determinant of the matrix,

$$\begin{pmatrix} 1 & 2 & 3 \\ 0 & -6 & -8 \\ 0 & 0 & -2 \end{pmatrix}$$

which equals 12. The cofactor matrix of  $A$  is

$$\begin{pmatrix} -2 & -2 & 6 \\ 4 & -2 & 0 \\ 2 & 8 & -6 \end{pmatrix}.$$

Each entry of  $A$  was replaced by its cofactor. Therefore, from the above theorem, the inverse of  $A$  should equal

$$\frac{1}{12} \begin{pmatrix} -2 & -2 & 6 \\ 4 & -2 & 0 \\ 2 & 8 & -6 \end{pmatrix}^T = \begin{pmatrix} -\frac{1}{6} & \frac{1}{3} & \frac{1}{2} \\ -\frac{1}{6} & -\frac{1}{6} & \frac{1}{3} \\ \frac{1}{2} & 0 & -\frac{1}{2} \end{pmatrix}.$$

This way of finding inverses is especially useful in the case where it is desired to find the inverse of a matrix whose entries are functions.

**Example 19.7.15** Suppose

$$A(t) = \begin{pmatrix} e^t & 0 & 0 \\ 0 & \cos t & \sin t \\ 0 & -\sin t & \cos t \end{pmatrix}$$

Find  $A(t)^{-1}$ .

First note  $\det(A(t)) = e^t$ . The cofactor matrix is

$$C(t) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & e^t \cos t & e^t \sin t \\ 0 & -e^t \sin t & e^t \cos t \end{pmatrix}$$

and so the inverse is

$$\frac{1}{e^t} \begin{pmatrix} 1 & 0 & 0 \\ 0 & e^t \cos t & e^t \sin t \\ 0 & -e^t \sin t & e^t \cos t \end{pmatrix}^T = \begin{pmatrix} e^{-t} & 0 & 0 \\ 0 & \cos t & -\sin t \\ 0 & \sin t & \cos t \end{pmatrix}.$$

This formula for the inverse also implies a famous procedure known as Cramer's rule. Cramer's rule gives a formula for the solutions,  $\mathbf{x}$ , to a system of equations,  $A\mathbf{x} = \mathbf{y}$ .

In case you are solving a system of equations,  $A\mathbf{x} = \mathbf{y}$  for  $\mathbf{x}$ , it follows that if  $A^{-1}$  exists,

$$\mathbf{x} = (A^{-1}A)\mathbf{x} = A^{-1}(A\mathbf{x}) = A^{-1}\mathbf{y}$$

thus solving the system. Now in the case that  $A^{-1}$  exists, there is a formula for  $A^{-1}$  given above. Using this formula,

$$x_i = \sum_{j=1}^n a_{ij}^{-1} y_j = \sum_{j=1}^n \frac{1}{\det(A)} \operatorname{cof}(A)_{ji} y_j.$$

By the formula for the expansion of a determinant along a column,

$$x_i = \frac{1}{\det(A)} \det \begin{pmatrix} * & \cdots & y_1 & \cdots & * \\ \vdots & & \vdots & & \vdots \\ * & \cdots & y_n & \cdots & * \end{pmatrix},$$

where here the  $i^{\text{th}}$  column of  $A$  is replaced with the column vector,  $(y_1, \dots, y_n)^T$ , and the determinant of this modified matrix is taken and divided by  $\det(A)$ . This formula is known as Cramer's rule.

**Procedure 19.7.16** Suppose  $A$  is an  $n \times n$  matrix and it is desired to solve the system  $A\mathbf{x} = \mathbf{y}$ ,  $\mathbf{y} = (y_1, \dots, y_n)^T$  for  $\mathbf{x} = (x_1, \dots, x_n)^T$ . Then Cramer's rule says

$$x_i = \frac{\det A_i}{\det A}$$

where  $A_i$  is obtained from  $A$  by replacing the  $i^{\text{th}}$  column of  $A$  with the column  $(y_1, \dots, y_n)^T$ .

**Definition 19.7.17** A submatrix of a matrix  $A$  is a rectangular array of numbers obtained by deleting some rows and columns of  $A$ . Let  $A$  be an  $m \times n$  matrix. The determinant rank of the matrix equals  $r$  where  $r$  is the largest number such that some  $r \times r$  submatrix of  $A$  has a non zero determinant. A given row,  $\mathbf{a}_s$  of a matrix,  $A$  is a linear combination of rows  $\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_r}$  if there are scalars,  $c_j$  such that  $\mathbf{a}_s = \sum_{j=1}^r c_j \mathbf{a}_{i_j}$ . The row rank of a matrix is the smallest number,  $r$  such that every row is a linear combination of some  $r$  rows. The column rank of a matrix is the smallest number,  $r$ , such that every column is a linear combination of some  $r$  columns.

The following theorem is proved in the section on the theory of the determinant. It says the row rank and column rank are both no larger than the determinant rank. In fact, it is easy to see they are all equal once the theorem has been proved.

**Theorem 19.7.18** *If  $A$  has determinant rank,  $r$ , then there exist  $r$  rows of the matrix such that every other row is a linear combination of these  $r$  rows. There also exist  $r$  columns such that every column is a linear combination of these  $r$  columns.*

The following theorem is of fundamental importance and ties together many of the ideas presented above. It is proved in the next section.

**Theorem 19.7.19** *Let  $A$  be an  $n \times n$  matrix. Then the following are equivalent.*

1.  $A$  is one to one.
2.  $A$  is onto.
3.  $\det(A) \neq 0$ .

It is a remarkable fact that for  $n \times n$  matrices, one to one is equivalent to onto. This amazing result is proved in the next section but is stated here.

**Corollary 19.7.20** *Let  $A$  be an  $n \times n$  matrix. Then  $A$  is one to one if and only if  $A$  is onto.*

## 19.8 Exercises

1. Find the determinants of the following matrices.

$$(a) \begin{pmatrix} 1 & 2 & 3 \\ 3 & 2 & 2 \\ 0 & 9 & 8 \end{pmatrix} \text{ (The answer is 31.)}$$

$$(b) \begin{pmatrix} 4 & 3 & 2 \\ 1 & 7 & 8 \\ 3 & -9 & 3 \end{pmatrix} \text{ (The answer is 375.)}$$

$$(c) \begin{pmatrix} 1 & 2 & 3 & 2 \\ 1 & 3 & 2 & 3 \\ 4 & 1 & 5 & 0 \\ 1 & 2 & 1 & 2 \end{pmatrix}, \text{ (The answer is } -2\text{.)}$$

2. A matrix is said to be orthogonal if  $A^T A = I$ . Thus the inverse of an orthogonal matrix is just its transpose. What are the possible values of  $\det(A)$  if  $A$  is an orthogonal matrix?
3. If  $A^{-1}$  exist, what is the relationship between  $\det(A)$  and  $\det(A^{-1})$ . Explain your answer.
4. Is it true that  $\det(A + B) = \det(A) + \det(B)$ ? If this is so, explain why it is so and if it is not so, give a counter example.
5. Let  $A$  be an  $r \times r$  matrix and suppose there are  $r - 1$  rows (columns) such that all rows (columns) are linear combinations of these  $r - 1$  rows (columns). Show  $\det(A) = 0$ .

6. †Theorem 19.7.18 showed that determinant rank is at least as large as column and row rank. Show they are all equal.
7. Show  $\det(aA) = a^n \det(A)$  where here  $A$  is an  $n \times n$  matrix and  $a$  is a scalar.
8. Suppose  $A$  is an upper triangular matrix. Show that  $A^{-1}$  exists if and only if all elements of the main diagonal are non zero. Is it true that  $A^{-1}$  will also be upper triangular? Explain. Is everything the same for lower triangular matrices?
9. Let  $A$  and  $B$  be two  $n \times n$  matrices.  $A \sim B$  ( $A$  is similar to  $B$ ) means there exists an invertible matrix,  $S$  such that  $A = S^{-1}BS$ . Show that if  $A \sim B$ , then  $B \sim A$ . Show also that  $A \sim A$  and that if  $A \sim B$  and  $B \sim C$ , then  $A \sim C$ .
10. In the context of Problem 9 show that if  $A \sim B$ , then  $\det(A) = \det(B)$ .
11. Let  $A$  be an  $n \times n$  matrix and let  $\mathbf{x}$  be a nonzero vector such that  $A\mathbf{x} = \lambda\mathbf{x}$  for some scalar,  $\lambda$ . When this occurs, the vector,  $\mathbf{x}$  is called an eigenvector and the scalar,  $\lambda$  is called an eigenvalue. It turns out that not every number is an eigenvalue. Only certain ones are. Why? **Hint:** Show that if  $A\mathbf{x} = \lambda\mathbf{x}$ , then  $(\lambda I - A)\mathbf{x} = \mathbf{0}$ . Explain why this shows that  $(\lambda I - A)$  is not one to one and not onto. Now use Theorem 19.7.19 to argue  $\det(\lambda I - A) = 0$ . What sort of equation is this? How many solutions does it have?
12. Suppose  $\det(\lambda I - A) = 0$ . Show using Theorem 19.7.19 there exists  $\mathbf{x} \neq \mathbf{0}$  such that  $(\lambda I - A)\mathbf{x} = \mathbf{0}$ .
13. Let  $F(t) = \det \begin{pmatrix} a(t) & b(t) \\ c(t) & d(t) \end{pmatrix}$ . Verify

$$F'(t) = \det \begin{pmatrix} a'(t) & b'(t) \\ c(t) & d(t) \end{pmatrix} + \det \begin{pmatrix} a(t) & b(t) \\ c'(t) & d'(t) \end{pmatrix}.$$

Now suppose

$$F(t) = \det \begin{pmatrix} a(t) & b(t) & c(t) \\ d(t) & e(t) & f(t) \\ g(t) & h(t) & i(t) \end{pmatrix}.$$

Use Laplace expansion and the first part to verify  $F'(t) =$

$$\det \begin{pmatrix} a'(t) & b'(t) & c'(t) \\ d(t) & e(t) & f(t) \\ g(t) & h(t) & i(t) \end{pmatrix} + \det \begin{pmatrix} a(t) & b(t) & c(t) \\ d'(t) & e'(t) & f'(t) \\ g(t) & h(t) & i(t) \end{pmatrix} + \det \begin{pmatrix} a(t) & b(t) & c(t) \\ d(t) & e(t) & f(t) \\ g'(t) & h'(t) & i'(t) \end{pmatrix}.$$

Conjecture a general result valid for  $n \times n$  matrices and explain why it will be true. Can a similar thing be done with the columns?

14. Use the formula for the inverse in terms of the cofactor matrix to find the inverse of the matrix,

$$A = \begin{pmatrix} e^t & 0 & 0 \\ 0 & e^t \cos t & e^t \sin t \\ 0 & e^t \cos t - e^t \sin t & e^t \cos t + e^t \sin t \end{pmatrix}.$$

15. Let  $A$  be an  $r \times r$  matrix and let  $B$  be an  $m \times m$  matrix such that  $r + m = n$ . Consider the following  $n \times n$  block matrix

$$C = \begin{pmatrix} A & 0 \\ D & B \end{pmatrix}.$$

where the  $D$  is an  $m \times r$  matrix, and the  $0$  is a  $r \times m$  matrix. Letting  $I_k$  denote the  $k \times k$  identity matrix, tell why

$$C = \begin{pmatrix} A & 0 \\ D & I_m \end{pmatrix} \begin{pmatrix} I_r & 0 \\ 0 & B \end{pmatrix}.$$

Now explain why  $\det(C) = \det(A) \det(B)$ . **Hint:** Part of this will require an explanation of why

$$\det \begin{pmatrix} A & 0 \\ D & I_m \end{pmatrix} = \det(A).$$

See Corollary 19.7.11.

## 19.9 The Mathematical Theory Of Determinants

It is easiest to give a different definition of the determinant which is clearly well defined and then prove the earlier one in terms of Laplace expansion. Let  $(i_1, \dots, i_n)$  be an ordered list of numbers from  $\{1, \dots, n\}$ . This means the order is important so  $(1, 2, 3)$  and  $(2, 1, 3)$  are different. There will be some repetition between this section and the previous section for the convenience of the reader.

The following Lemma will be essential in the definition of the determinant.

**Lemma 19.9.1** *There exists a unique function,  $\text{sgn}_n$  which maps each list of numbers from  $\{1, \dots, n\}$  to one of the three numbers,  $0, 1$ , or  $-1$  which also has the following properties.*

$$\text{sgn}_n(1, \dots, n) = 1 \quad (19.40)$$

$$\text{sgn}_n(i_1, \dots, p, \dots, q, \dots, i_n) = -\text{sgn}_n(i_1, \dots, q, \dots, p, \dots, i_n) \quad (19.41)$$

*In words, the second property states that if two of the numbers are switched, the value of the function is multiplied by  $-1$ . Also, in the case where  $n > 1$  and  $\{i_1, \dots, i_n\} = \{1, \dots, n\}$  so that every number from  $\{1, \dots, n\}$  appears in the ordered list,  $(i_1, \dots, i_n)$ ,*

$$\begin{aligned} \text{sgn}_n(i_1, \dots, i_{\theta-1}, n, i_{\theta+1}, \dots, i_n) &\equiv \\ (-1)^{n-\theta} \text{sgn}_{n-1}(i_1, \dots, i_{\theta-1}, i_{\theta+1}, \dots, i_n) &\quad (19.42) \end{aligned}$$

*where  $n = i_\theta$  in the ordered list,  $(i_1, \dots, i_n)$ .*

**Proof:** To begin with, it is necessary to show the existence of such a function. This is clearly true if  $n = 1$ . Define  $\text{sgn}_1(1) \equiv 1$  and observe that it works. No switching is possible. In the case where  $n = 2$ , it is also clearly true. Let  $\text{sgn}_2(1, 2) = 1$  and  $\text{sgn}_2(2, 1) = 0$  while  $\text{sgn}_2(2, 2) = \text{sgn}_2(1, 1) = 0$  and verify it works. Assuming such a function exists for  $n$ ,  $\text{sgn}_{n+1}$  will be defined in terms of  $\text{sgn}_n$ . If there are any repeated numbers in  $(i_1, \dots, i_{n+1})$ ,  $\text{sgn}_{n+1}(i_1, \dots, i_{n+1}) \equiv 0$ . If there are no repeats, then  $n+1$  appears somewhere in the ordered list. Let  $\theta$  be the position of the number  $n+1$  in the list. Thus, the list is of the form  $(i_1, \dots, i_{\theta-1}, n+1, i_{\theta+1}, \dots, i_{n+1})$ . From (19.42) it must be that

$$\text{sgn}_{n+1}(i_1, \dots, i_{\theta-1}, n+1, i_{\theta+1}, \dots, i_{n+1}) \equiv$$

$$(-1)^{n+1-\theta} \operatorname{sgn}_n(i_1, \dots, i_{\theta-1}, i_{\theta+1}, \dots, i_{n+1}).$$

It is necessary to verify this satisfies (19.40) and (19.41) with  $n$  replaced with  $n+1$ . The first of these is obviously true because

$$\operatorname{sgn}_{n+1}(i_1, \dots, n, n+1) \equiv (-1)^{n+1-(n+1)} \operatorname{sgn}_n(i_1, \dots, n) = 1.$$

If there are repeated numbers in  $(i_1, \dots, i_{n+1})$ , then it is obvious (19.41) holds because both sides would equal zero from the above definition. It remains to verify (19.41) in the case where there are no numbers repeated in  $(i_1, \dots, i_{n+1})$ . Consider

$$\operatorname{sgn}_{n+1}(i_1, \dots, \overset{r}{p}, \dots, \overset{s}{q}, \dots, i_{n+1}),$$

where the  $r$  above the  $p$  indicates the number,  $p$  is in the  $r^{\text{th}}$  position and the  $s$  above the  $q$  indicates that the number,  $q$  is in the  $s^{\text{th}}$  position. Suppose first that  $r < \theta < s$ . Then

$$\begin{aligned} \operatorname{sgn}_{n+1}(i_1, \dots, \overset{r}{p}, \dots, \overset{\theta}{n+1}, \dots, \overset{s}{q}, \dots, i_{n+1}) &\equiv \\ (-1)^{n+1-\theta} \operatorname{sgn}_n(i_1, \dots, \overset{r}{p}, \dots, \overset{s-1}{q}, \dots, i_{n+1}) \end{aligned}$$

while

$$\begin{aligned} \operatorname{sgn}_{n+1}(i_1, \dots, \overset{r}{q}, \dots, \overset{\theta}{n+1}, \dots, \overset{s}{p}, \dots, i_{n+1}) &= \\ (-1)^{n+1-\theta} \operatorname{sgn}_n(i_1, \dots, \overset{r}{q}, \dots, \overset{s-1}{p}, \dots, i_{n+1}) \end{aligned}$$

and so, by induction, a switch of  $p$  and  $q$  introduces a minus sign in the result. Similarly, if  $\theta > s$  or if  $\theta < r$  it also follows that (19.41) holds. The interesting case is when  $\theta = r$  or  $\theta = s$ . Consider the case where  $\theta = r$  and note the other case is entirely similar.

$$\begin{aligned} \operatorname{sgn}_{n+1}(i_1, \dots, \overset{r}{n+1}, \dots, \overset{s}{q}, \dots, i_{n+1}) &= \\ (-1)^{n+1-r} \operatorname{sgn}_n(i_1, \dots, \overset{s-1}{q}, \dots, i_{n+1}) \end{aligned} \tag{19.43}$$

while

$$\begin{aligned} \operatorname{sgn}_{n+1}(i_1, \dots, \overset{r}{q}, \dots, \overset{s}{n+1}, \dots, i_{n+1}) &= \\ (-1)^{n+1-s} \operatorname{sgn}_n(i_1, \dots, \overset{r}{q}, \dots, i_{n+1}). \end{aligned} \tag{19.44}$$

By making  $s-1-r$  switches, move the  $q$  which is in the  $s-1^{\text{th}}$  position in (19.43) to the  $r^{\text{th}}$  position in (19.44). By induction, each of these switches introduces a factor of  $-1$  and so

$$\operatorname{sgn}_n(i_1, \dots, \overset{s-1}{q}, \dots, i_{n+1}) = (-1)^{s-1-r} \operatorname{sgn}_n(i_1, \dots, \overset{r}{q}, \dots, i_{n+1}).$$

Therefore,

$$\begin{aligned} \operatorname{sgn}_{n+1}(i_1, \dots, \overset{r}{n+1}, \dots, \overset{s}{q}, \dots, i_{n+1}) &= (-1)^{n+1-r} \operatorname{sgn}_n(i_1, \dots, \overset{s-1}{q}, \dots, i_{n+1}) \\ &= (-1)^{n+1-r} (-1)^{s-1-r} \operatorname{sgn}_n(i_1, \dots, \overset{r}{q}, \dots, i_{n+1}) \\ &= (-1)^{n+s} \operatorname{sgn}_n(i_1, \dots, \overset{r}{q}, \dots, i_{n+1}) = (-1)^{2s-1} (-1)^{n+1-s} \operatorname{sgn}_n(i_1, \dots, \overset{r}{q}, \dots, i_{n+1}) \end{aligned}$$



$$= -\operatorname{sgn}_{n+1} \left( i_1, \dots, \overset{r}{q}, \dots, n + \overset{s}{1}, \dots, i_{n+1} \right).$$

This proves the existence of the desired function.

To see this function is unique, note that you can obtain any ordered list of distinct numbers from a sequence of switches. If there exist two functions,  $f$  and  $g$  both satisfying (19.40) and (19.41), you could start with  $f(1, \dots, n) = g(1, \dots, n)$  and applying the same sequence of switches, eventually arrive at  $f(i_1, \dots, i_n) = g(i_1, \dots, i_n)$ . If any numbers are repeated, then (19.41) gives both functions are equal to zero for that ordered list. This proves the lemma.

In what follows  $\operatorname{sgn}$  will often be used rather than  $\operatorname{sgn}_n$  because the context supplies the appropriate  $n$ .

**Definition 19.9.2** Let  $f$  be a real valued function which has the set of ordered lists of numbers from  $\{1, \dots, n\}$  as its domain. Define

$$\sum_{(k_1, \dots, k_n)} f(k_1 \cdots k_n)$$

to be the sum of all the  $f(k_1 \cdots k_n)$  for all possible choices of ordered lists  $(k_1, \dots, k_n)$  of numbers of  $\{1, \dots, n\}$ . For example,

$$\sum_{(k_1, k_2)} f(k_1, k_2) = f(1, 2) + f(2, 1) + f(1, 1) + f(2, 2).$$

**Definition 19.9.3** Let  $(a_{ij}) = A$  denote an  $n \times n$  matrix. The determinant of  $A$ , denoted by  $\det(A)$  is defined by

$$\det(A) \equiv \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) a_{1k_1} \cdots a_{nk_n}$$

where the sum is taken over all ordered lists of numbers from  $\{1, \dots, n\}$ . Note it suffices to take the sum over only those ordered lists in which there are no repeats because if there are,  $\operatorname{sgn}(k_1, \dots, k_n) = 0$  and so that term contributes 0 to the sum.

Let  $A$  be an  $n \times n$  matrix,  $A = (a_{ij})$  and let  $(r_1, \dots, r_n)$  denote an ordered list of  $n$  numbers from  $\{1, \dots, n\}$ . Let  $A(r_1, \dots, r_n)$  denote the matrix whose  $k^{\text{th}}$  row is the  $r_k$  row of the matrix,  $A$ . Thus

$$\det(A(r_1, \dots, r_n)) = \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) a_{r_1 k_1} \cdots a_{r_n k_n} \quad (19.45)$$

and

$$A(1, \dots, n) = A.$$

**Proposition 19.9.4** Let

$$(r_1, \dots, r_n)$$

be an ordered list of numbers from  $\{1, \dots, n\}$ . Then

$$\operatorname{sgn}(r_1, \dots, r_n) \det(A)$$

$$= \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) a_{r_1 k_1} \cdots a_{r_n k_n} \quad (19.46)$$

$$= \det(A(r_1, \dots, r_n)). \quad (19.47)$$

**Proof:** Let  $(1, \dots, n) = (1, \dots, r, \dots, s, \dots, n)$  so  $r < s$ .

$$\det(A(1, \dots, r, \dots, s, \dots, n)) = \quad (19.48)$$

$$\sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_r, \dots, k_s, \dots, k_n) a_{1k_1} \cdots a_{rk_r} \cdots a_{sk_s} \cdots a_{nk_n},$$

and renaming the variables, calling  $k_s, k_r$  and  $k_r, k_s$ , this equals

$$\begin{aligned} &= \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_s, \dots, k_r, \dots, k_n) a_{1k_1} \cdots a_{rk_s} \cdots a_{sk_r} \cdots a_{nk_n} \\ &= \sum_{(k_1, \dots, k_n)} -\operatorname{sgn}\left(k_1, \dots, \overbrace{k_r, \dots, k_s}^{\text{These got switched}}, \dots, k_n\right) a_{1k_1} \cdots a_{sk_r} \cdots a_{rk_s} \cdots a_{nk_n} \\ &= -\det(A(1, \dots, s, \dots, r, \dots, n)). \end{aligned} \quad (19.49)$$

Consequently,

$$\begin{aligned} \det(A(1, \dots, s, \dots, r, \dots, n)) &= \\ &= -\det(A(1, \dots, r, \dots, s, \dots, n)) = -\det(A) \end{aligned}$$

Now letting  $A(1, \dots, s, \dots, r, \dots, n)$  play the role of  $A$ , and continuing in this way, switching pairs of numbers,

$$\det(A(r_1, \dots, r_n)) = (-1)^p \det(A)$$

where it took  $p$  switches to obtain  $(r_1, \dots, r_n)$  from  $(1, \dots, n)$ . By Lemma 19.9.1, this implies

$$\det(A(r_1, \dots, r_n)) = (-1)^p \det(A) = \operatorname{sgn}(r_1, \dots, r_n) \det(A)$$

and proves the proposition in the case when there are no repeated numbers in the ordered list,  $(r_1, \dots, r_n)$ . However, if there is a repeat, say the  $r^{\text{th}}$  row equals the  $s^{\text{th}}$  row, then the reasoning of (19.48)–(19.49) shows that  $A(r_1, \dots, r_n) = 0$  and also  $\operatorname{sgn}(r_1, \dots, r_n) = 0$  so the formula holds in this case also.

**Observation 19.9.5** *There are  $n!$  ordered lists of distinct numbers from  $\{1, \dots, n\}$ .*

To see this, consider  $n$  slots placed in order. There are  $n$  choices for the first slot. For each of these choices, there are  $n - 1$  choices for the second. Thus there are  $n(n - 1)$  ways to fill the first two slots. Then for each of these ways there are  $n - 2$  choices left for the third slot. Continuing this way, there are  $n!$  ordered lists of distinct numbers from  $\{1, \dots, n\}$  as stated in the observation.

With the above, it is possible to give a more symmetric description of the determinant from which it will follow that  $\det(A) = \det(A^T)$ .

**Corollary 19.9.6** *The following formula for  $\det(A)$  is valid.*

$$\det(A) = \frac{1}{n!} \sum_{(r_1, \dots, r_n)} \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(r_1, \dots, r_n) \operatorname{sgn}(k_1, \dots, k_n) a_{r_1 k_1} \cdots a_{r_n k_n}. \quad (19.50)$$

And also  $\det(A^T) = \det(A)$  where  $A^T$  is the transpose of  $A$ . (Recall that for  $A^T = (a_{ij}^T)$ ,  $a_{ij}^T = a_{ji}$ .)

**Proof:** From Proposition 19.9.4, if the  $r_i$  are distinct,

$$\det(A) = \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(r_1, \dots, r_n) \operatorname{sgn}(k_1, \dots, k_n) a_{r_1 k_1} \cdots a_{r_n k_n}.$$

Summing over all ordered lists,  $(r_1, \dots, r_n)$  where the  $r_i$  are distinct, (If the  $r_i$  are not distinct,  $\operatorname{sgn}(r_1, \dots, r_n) = 0$  and so there is no contribution to the sum.)

$$n! \det(A) =$$

$$\sum_{(r_1, \dots, r_n)} \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(r_1, \dots, r_n) \operatorname{sgn}(k_1, \dots, k_n) a_{r_1 k_1} \cdots a_{r_n k_n}.$$

This proves the corollary since the formula gives the same number for  $A$  as it does for  $A^T$ .

**Corollary 19.9.7** *If two rows or two columns in an  $n \times n$  matrix,  $A$ , are switched, the determinant of the resulting matrix equals  $(-1)$  times the determinant of the original matrix. If  $A$  is an  $n \times n$  matrix in which two rows are equal or two columns are equal then  $\det(A) = 0$ . Suppose the  $i^{\text{th}}$  row of  $A$  equals  $(xa_1 + yb_1, \dots, xa_n + yb_n)$ . Then*

$$\det(A) = x \det(A_1) + y \det(A_2)$$

where the  $i^{\text{th}}$  row of  $A_1$  is  $(a_1, \dots, a_n)$  and the  $i^{\text{th}}$  row of  $A_2$  is  $(b_1, \dots, b_n)$ , all other rows of  $A_1$  and  $A_2$  coinciding with those of  $A$ . In other words,  $\det$  is a linear function of each row  $A$ . The same is true with the word “row” replaced with the word “column”.

**Proof:** By Proposition 19.9.4 when two rows are switched, the determinant of the resulting matrix is  $(-1)$  times the determinant of the original matrix. By Corollary 19.9.6 the same holds for columns because the columns of the matrix equal the rows of the transposed matrix. Thus if  $A_1$  is the matrix obtained from  $A$  by switching two columns,

$$\det(A) = \det(A^T) = -\det(A_1^T) = -\det(A_1).$$

If  $A$  has two equal columns or two equal rows, then switching them results in the same matrix. Therefore,  $\det(A) = -\det(A)$  and so  $\det(A) = 0$ .

It remains to verify the last assertion.

$$\begin{aligned} \det(A) &\equiv \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) a_{1k_1} \cdots (xa_{ik_i} + yb_{ik_i}) \cdots a_{nk_n} \\ &= x \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) a_{1k_1} \cdots a_{ik_i} \cdots a_{nk_n} \\ &\quad + y \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) a_{1k_1} \cdots b_{ik_i} \cdots a_{nk_n} \\ &\equiv x \det(A_1) + y \det(A_2). \end{aligned}$$

The same is true of columns because  $\det(A^T) = \det(A)$  and the rows of  $A^T$  are the columns of  $A$ .

**Definition 19.9.8** *We say a vector,  $\mathbf{w}$ , is a linear combination of the vectors  $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$  if there exists scalars,  $c_1, \dots, c_r$  such that  $\mathbf{w} = \sum_{k=1}^r c_k \mathbf{v}_k$ . This is the same as saying  $\mathbf{w} \in \operatorname{span}\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$ .*

The following corollary is also of great use.

**Corollary 19.9.9** Suppose  $A$  is an  $n \times n$  matrix and some column (row) is a linear combination of  $r$  other columns (rows). Then  $\det(A) = 0$ .

**Proof:** Let  $A = (\mathbf{a}_1 \cdots \mathbf{a}_n)$  be the columns of  $A$  and suppose the condition that one column is a linear combination of  $r$  of the others is satisfied. Then by using Corollary 19.9.7 we may rearrange the columns to have the  $n^{\text{th}}$  column a linear combination of the first  $r$  columns. Thus  $\mathbf{a}_n = \sum_{k=1}^r c_k \mathbf{a}_k$  and so

$$\det(A) = \det(\mathbf{a}_1 \cdots \mathbf{a}_r \cdots \mathbf{a}_{n-1} \sum_{k=1}^r c_k \mathbf{a}_k).$$

By Corollary 19.9.7

$$\det(A) = \sum_{k=1}^r c_k \det(\mathbf{a}_1 \cdots \mathbf{a}_r \cdots \mathbf{a}_{n-1} \mathbf{a}_k) = 0.$$

The case for rows follows from the fact that  $\det(A) = \det(A^T)$ . This proves the corollary.

Recall the following definition of matrix multiplication.

**Definition 19.9.10** If  $A$  and  $B$  are  $n \times n$  matrices,  $A = (a_{ij})$  and  $B = (b_{ij})$ ,  $AB = (c_{ij})$  where

$$c_{ij} \equiv \sum_{k=1}^n a_{ik} b_{kj}.$$

One of the most important rules about determinants is that the determinant of a product equals the product of the determinants.

**Theorem 19.9.11** Let  $A$  and  $B$  be  $n \times n$  matrices. Then

$$\det(AB) = \det(A) \det(B).$$

**Proof:** Let  $c_{ij}$  be the  $ij^{\text{th}}$  entry of  $AB$ . Then by Proposition 19.9.4,

$$\begin{aligned} \det(AB) &= \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) c_{1k_1} \cdots c_{nk_n} \\ &= \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) \left( \sum_{r_1} a_{1r_1} b_{r_1 k_1} \right) \cdots \left( \sum_{r_n} a_{nr_n} b_{r_n k_n} \right) \\ &= \sum_{(r_1, \dots, r_n)} \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) b_{r_1 k_1} \cdots b_{r_n k_n} (a_{1r_1} \cdots a_{nr_n}) \\ &= \sum_{(r_1, \dots, r_n)} \operatorname{sgn}(r_1 \cdots r_n) a_{1r_1} \cdots a_{nr_n} \det(B) = \det(A) \det(B). \end{aligned}$$

This proves the theorem.

**Lemma 19.9.12** Suppose a matrix is of the form

$$M = \begin{pmatrix} A & * \\ \mathbf{0} & a \end{pmatrix} \quad (19.51)$$

or

$$M = \begin{pmatrix} A & \mathbf{0} \\ * & a \end{pmatrix} \quad (19.52)$$

where  $a$  is a number and  $A$  is an  $(n-1) \times (n-1)$  matrix and  $*$  denotes either a column or a row having length  $n-1$  and the  $\mathbf{0}$  denotes either a column or a row of length  $n-1$  consisting entirely of zeros. Then

$$\det(M) = a \det(A).$$

**Proof:** Denote  $M$  by  $(m_{ij})$ . Thus in the first case,  $m_{nn} = a$  and  $m_{ni} = 0$  if  $i \neq n$  while in the second case,  $m_{nn} = a$  and  $m_{in} = 0$  if  $i \neq n$ . From the definition of the determinant,

$$\det(M) \equiv \sum_{(k_1, \dots, k_n)} \operatorname{sgn}_n(k_1, \dots, k_n) m_{1k_1} \cdots m_{nk_n}$$

Letting  $\theta$  denote the position of  $n$  in the ordered list,  $(k_1, \dots, k_n)$  then using the earlier conventions used to prove Lemma 19.9.1,  $\det(M)$  equals

$$\sum_{(k_1, \dots, k_n)} (-1)^{n-\theta} \operatorname{sgn}_{n-1} \left( k_1, \dots, k_{\theta-1}, k_{\theta+1}, \dots, k_n \right) m_{1k_1} \cdots m_{nk_n}$$

Now suppose (19.52). Then if  $k_n \neq n$ , the term involving  $m_{nk_n}$  in the above expression equals zero. Therefore, the only terms which survive are those for which  $\theta = n$  or in other words, those for which  $k_n = n$ . Therefore, the above expression reduces to

$$a \sum_{(k_1, \dots, k_{n-1})} \operatorname{sgn}_{n-1}(k_1, \dots, k_{n-1}) m_{1k_1} \cdots m_{(n-1)k_{n-1}} = a \det(A).$$

To get the assertion in the situation of (19.51) use Corollary 19.9.6 and (19.52) to write

$$\det(M) = \det(M^T) = \det \left( \begin{pmatrix} A^T & \mathbf{0} \\ * & a \end{pmatrix} \right) = a \det(A^T) = a \det(A).$$

This proves the lemma.

In terms of the theory of determinants, arguably the most important idea is that of Laplace expansion along a row or a column. This will follow from the above definition of a determinant.

**Definition 19.9.13** Let  $A = (a_{ij})$  be an  $n \times n$  matrix. Then a new matrix called the cofactor matrix,  $\operatorname{cof}(A)$  is defined by  $\operatorname{cof}(A) = (c_{ij})$  where to obtain  $c_{ij}$  delete the  $i^{\text{th}}$  row and the  $j^{\text{th}}$  column of  $A$ , take the determinant of the  $(n-1) \times (n-1)$  matrix which results, (This is called the  $ij^{\text{th}}$  minor of  $A$ .) and then multiply this number by  $(-1)^{i+j}$ . To make the formulas easier to remember,  $\operatorname{cof}(A)_{ij}$  will denote the  $ij^{\text{th}}$  entry of the cofactor matrix.

The following is the main result. Earlier this was given as a definition and the outrageous totally unjustified assertion was made that the same number would be obtained by expanding the determinant along any row or column. The following theorem proves this assertion.

**Theorem 19.9.14** Let  $A$  be an  $n \times n$  matrix where  $n \geq 2$ . Then

$$\det(A) = \sum_{j=1}^n a_{ij} \operatorname{cof}(A)_{ij} = \sum_{i=1}^n a_{ij} \operatorname{cof}(A)_{ij}. \quad (19.53)$$

The first formula consists of expanding the determinant along the  $i^{\text{th}}$  row and the second expands the determinant along the  $j^{\text{th}}$  column.

**Proof:** Let  $(a_{i1}, \dots, a_{in})$  be the  $i^{th}$  row of  $A$ . Let  $B_j$  be the matrix obtained from  $A$  by leaving every row the same except the  $i^{th}$  row which in  $B_j$  equals  $(0, \dots, 0, a_{ij}, 0, \dots, 0)$ . Then by Corollary 19.9.7,

$$\det(A) = \sum_{j=1}^n \det(B_j)$$

Denote by  $A^{ij}$  the  $(n-1) \times (n-1)$  matrix obtained by deleting the  $i^{th}$  row and the  $j^{th}$  column of  $A$ . Thus  $\text{cof}(A)_{ij} \equiv (-1)^{i+j} \det(A^{ij})$ . At this point, recall that from Proposition 19.9.4, when two rows or two columns in a matrix,  $M$ , are switched, this results in multiplying the determinant of the old matrix by  $-1$  to get the determinant of the new matrix. Therefore, by Lemma 19.9.12,

$$\begin{aligned} \det(B_j) &= (-1)^{n-j} (-1)^{n-i} \det \left( \begin{pmatrix} A^{ij} & * \\ \mathbf{0} & a_{ij} \end{pmatrix} \right) \\ &= (-1)^{i+j} \det \left( \begin{pmatrix} A^{ij} & * \\ \mathbf{0} & a_{ij} \end{pmatrix} \right) = a_{ij} \text{cof}(A)_{ij}. \end{aligned}$$

Therefore,

$$\det(A) = \sum_{j=1}^n a_{ij} \text{cof}(A)_{ij}$$

which is the formula for expanding  $\det(A)$  along the  $i^{th}$  row. Also,

$$\begin{aligned} \det(A) &= \det(A^T) = \sum_{j=1}^n a_{ij}^T \text{cof}(A^T)_{ij} \\ &= \sum_{j=1}^n a_{ji} \text{cof}(A)_{ji} \end{aligned}$$

which is the formula for expanding  $\det(A)$  along the  $i^{th}$  column. This proves the theorem.

Note that this gives an easy way to write a formula for the inverse of an  $n \times n$  matrix. Recall the definition of the inverse of a matrix in Definition 19.1.13 on Page 417.

**Theorem 19.9.15**  $A^{-1}$  exists if and only if  $\det(A) \neq 0$ . If  $\det(A) \neq 0$ , then  $A^{-1} = (a_{ij}^{-1})$  where

$$a_{ij}^{-1} = \det(A)^{-1} \text{cof}(A)_{ji}$$

for  $\text{cof}(A)_{ij}$  the  $ij^{th}$  cofactor of  $A$ .

**Proof:** By Theorem 19.9.14 and letting  $(a_{ir}) = A$ , if  $\det(A) \neq 0$ ,

$$\sum_{i=1}^n a_{ir} \text{cof}(A)_{ir} \det(A)^{-1} = \det(A) \det(A)^{-1} = 1.$$

Now consider

$$\sum_{i=1}^n a_{ir} \text{cof}(A)_{ik} \det(A)^{-1}$$

when  $k \neq r$ . Replace the  $k^{th}$  column with the  $r^{th}$  column to obtain a matrix,  $B_k$  whose determinant equals zero by Corollary 19.9.7. However, expanding this matrix along the  $k^{th}$  column yields

$$0 = \det(B_k) \det(A)^{-1} = \sum_{i=1}^n a_{ir} \text{cof}(A)_{ik} \det(A)^{-1}$$

Summarizing,

$$\sum_{i=1}^n a_{ir} \operatorname{cof}(A)_{ik} \det(A)^{-1} = \delta_{rk}.$$

Using the other formula in Theorem 19.9.14, and similar reasoning,

$$\sum_{j=1}^n a_{rj} \operatorname{cof}(A)_{kj} \det(A)^{-1} = \delta_{rk}$$

This proves that if  $\det(A) \neq 0$ , then  $A^{-1}$  exists with  $A^{-1} = (a_{ij}^{-1})$ , where

$$a_{ij}^{-1} = \operatorname{cof}(A)_{ji} \det(A)^{-1}.$$

Now suppose  $A^{-1}$  exists. Then by Theorem 19.9.11,

$$1 = \det(I) = \det(AA^{-1}) = \det(A) \det(A^{-1})$$

so  $\det(A) \neq 0$ . This proves the theorem.

The next corollary points out that if an  $n \times n$  matrix,  $A$  has a right or a left inverse, then it has an inverse.

**Corollary 19.9.16** *Let  $A$  be an  $n \times n$  matrix and suppose there exists an  $n \times n$  matrix,  $B$  such that  $BA = I$ . Then  $A^{-1}$  exists and  $A^{-1} = B$ . Also, if there exists  $C$  an  $n \times n$  matrix such that  $AC = I$ , then  $A^{-1}$  exists and  $A^{-1} = C$ .*

**Proof:** Since  $BA = I$ , Theorem 19.9.11 implies

$$\det B \det A = 1$$

and so  $\det A \neq 0$ . Therefore from Theorem 19.9.15,  $A^{-1}$  exists. Therefore,

$$A^{-1} = (BA) A^{-1} = B(AA^{-1}) = BI = B.$$

The case where  $CA = I$  is handled similarly.

The conclusion of this corollary is that left inverses, right inverses and inverses are all the same in the context of  $n \times n$  matrices.

Theorem 19.9.15 says that to find the inverse, take the transpose of the cofactor matrix and divide by the determinant. The transpose of the cofactor matrix is called the adjugate or sometimes the classical adjoint of the matrix  $A$ . It is an abomination to call it the adjoint although you do sometimes see it referred to in this way. In words,  $A^{-1}$  is equal to one over the determinant of  $A$  times the adjugate matrix of  $A$ .

In case you are solving a system of equations,  $A\mathbf{x} = \mathbf{y}$  for  $\mathbf{x}$ , it follows that if  $A^{-1}$  exists,

$$\mathbf{x} = (A^{-1}A) \mathbf{x} = A^{-1}(A\mathbf{x}) = A^{-1}\mathbf{y}$$

thus solving the system. Now in the case that  $A^{-1}$  exists, there is a formula for  $A^{-1}$  given above. Using this formula,

$$x_i = \sum_{j=1}^n a_{ij}^{-1} y_j = \sum_{j=1}^n \frac{1}{\det(A)} \operatorname{cof}(A)_{ji} y_j.$$

By the formula for the expansion of a determinant along a column,

$$x_i = \frac{1}{\det(A)} \det \begin{pmatrix} * & \cdots & y_1 & \cdots & * \\ \vdots & & \vdots & & \vdots \\ * & \cdots & y_n & \cdots & * \end{pmatrix},$$

where here the  $i^{\text{th}}$  column of  $A$  is replaced with the column vector,  $(y_1 \cdots y_n)^T$ , and the determinant of this modified matrix is taken and divided by  $\det(A)$ . This formula is known as Cramer's rule.

**Definition 19.9.17** A matrix  $M$ , is upper triangular if  $M_{ij} = 0$  whenever  $i > j$ . Thus such a matrix equals zero below the main diagonal, the entries of the form  $M_{ii}$  as shown.

$$\begin{pmatrix} * & * & \cdots & * \\ 0 & * & \ddots & \vdots \\ \vdots & \ddots & \ddots & * \\ 0 & \cdots & 0 & * \end{pmatrix}$$

A lower triangular matrix is defined similarly as a matrix for which all entries above the main diagonal are equal to zero.

With this definition, here is a simple corollary of Theorem 19.9.14.

**Corollary 19.9.18** Let  $M$  be an upper (lower) triangular matrix. Then  $\det(M)$  is obtained by taking the product of the entries on the main diagonal.

Recall the following definition of rank presented earlier.

**Definition 19.9.19** A submatrix of a matrix  $A$  is a rectangular array of numbers obtained by deleting some rows and columns of  $A$ . Let  $A$  be an  $m \times n$  matrix. The determinant rank of the matrix equals  $r$  where  $r$  is the largest number such that some  $r \times r$  submatrix of  $A$  has a non zero determinant. A given row,  $\mathbf{a}_s$  of a matrix,  $A$  is a linear combination of rows  $\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_r}$  if there are scalars,  $c_j$  such that  $\mathbf{a}_s = \sum_{j=1}^r c_j \mathbf{a}_{i_j}$ . The row rank of a matrix is the smallest number,  $r$  such that every row is a linear combination of some  $r$  rows. The column rank of a matrix is the smallest number,  $r$ , such that every column is a linear combination of some  $r$  columns.

**Theorem 19.9.20** The determinant rank of  $A$  coincides with the row rank.

**Proof:** Suppose the determinant rank of  $A = (a_{ij})$  equals  $r$ . If rows and columns are interchanged, the determinant rank of the modified matrix is unchanged. Thus rows and columns can be interchanged to produce an  $r \times r$  matrix in the upper left corner of the matrix which has non zero determinant. Now consider the  $r+1 \times r+1$  matrix,  $M$ ,

$$\begin{pmatrix} a_{11} & \cdots & a_{1r} & a_{1p} \\ \vdots & & \vdots & \vdots \\ a_{r1} & \cdots & a_{rr} & a_{rp} \\ a_{l1} & \cdots & a_{lr} & a_{lp} \end{pmatrix}$$

where  $C$  will denote the  $r \times r$  matrix in the upper left corner which has non zero determinant. I claim  $\det(M) = 0$ .

There are two cases to consider in verifying this claim. First, suppose  $p > r$ . Then the claim follows from the assumption that  $A$  has determinant rank  $r$ . On the other hand, if  $p < r$ , then the determinant is zero because there are two identical columns. Expand the determinant along the last column and divide by  $\det(C)$  to obtain

$$a_{lp} = - \sum_{i=1}^r \frac{\text{cof}(M)_{ip}}{\det(C)} a_{ip}.$$



Now note that  $\text{cof}(M)_{ip}$  does not depend on  $p$ . Therefore the above sum is of the form

$$a_{lp} = \sum_{i=1}^r m_i a_{ip}$$

which shows the  $l^{\text{th}}$  row is a linear combination of the first  $r$  rows of  $A$ . Since  $l$  is arbitrary, this proves that every row is a linear combination of the first  $r$  rows.

Now suppose for some  $k < r$ , there exist  $k$  rows such that every row is a linear combination of these  $k$  rows. Then the matrix,  $C$  in the upper left corner must have determinant equal to zero because this matrix would be of the form

$$\begin{pmatrix} \sum_{i=1}^k c_i^1 \mathbf{v}_i^T \\ \vdots \\ \sum_{i=1}^k c_i^r \mathbf{v}_i^T \end{pmatrix}$$

where the rows are  $\{\mathbf{v}_i^T\}_{i=1}^k$ . By Corollary 19.9.7 on Page 451 the determinant of the above matrix is of the form

$$\sum_{i_1, \dots, i_k} c_{i_1}^1 c_{i_2}^2 \cdots c_{i_k}^k \det \begin{pmatrix} \mathbf{v}_{i_1}^T \\ \vdots \\ \mathbf{v}_{i_k}^T \end{pmatrix}$$

where each matrix in the above sum has  $r$  rows. But there are only  $k$  choices for vectors to fill these rows and so there must be repeated rows resulting in the above sum equaling zero as claimed. Therefore,  $k \geq r$  after all and so the row rank of  $A$  equals  $r$ .

**Corollary 19.9.21** *The column rank of  $A$  coincides with the determinant rank of  $A$ .*

**Proof:** This follows from the above theorem by considering  $A^T$ . The rows of  $A^T$  are the columns of  $A$  and the determinant rank of  $A^T$  and  $A$  are the same by Corollary 19.9.6 on Page 450.

The following theorem is of fundamental importance and ties together many of the ideas presented above.

**Theorem 19.9.22** *Let  $A$  be an  $n \times n$  matrix. Then the following are equivalent.*

1.  $\det(A) = 0$ .
2.  $A, A^T$  are not one to one.
3.  $A$  is not onto.

**Proof:** Suppose  $\det(A) = 0$ . Then the determinant rank of  $A = r < n$ . Therefore, there exist  $r$  columns such that every other column is a linear combination of these columns by Corollary 19.9.21. In particular, it follows that for some  $m$ , the  $m^{\text{th}}$  column is a linear combination of all the others. Thus letting  $A = (\mathbf{a}_1 \cdots \mathbf{a}_m \cdots \mathbf{a}_n)$  where the columns are denoted by  $\mathbf{a}_i$ , there exists scalars,  $\alpha_i$  such that

$$\mathbf{a}_m = \sum_{k \neq m} \alpha_k \mathbf{a}_k.$$

Now consider the column vector,  $\mathbf{x} \equiv (\alpha_1 \cdots -1 \cdots \alpha_n)^T$ . Then

$$A\mathbf{x} = -\mathbf{a}_m + \sum_{k \neq m} \alpha_k \mathbf{a}_k = \mathbf{0}.$$

Since also  $A\mathbf{0} = \mathbf{0}$ , it follows  $A$  is not one to one. Similarly,  $A^T$  is not one to one by the same argument applied to  $A^T$ . This verifies that 1.) implies 2.).

Now suppose 2.). Then since  $A^T$  is not one to one, it follows there exists  $\mathbf{x} \neq \mathbf{0}$  such that

$$A^T \mathbf{x} = \mathbf{0}.$$

Taking the transpose of both sides yields

$$\mathbf{x}^T A = \mathbf{0}$$

where the  $\mathbf{0}$  is a  $1 \times n$  matrix or row vector. Now if  $A\mathbf{y} = \mathbf{x}$ , then

$$|\mathbf{x}|^2 = \mathbf{x}^T (A\mathbf{y}) = (\mathbf{x}^T A) \mathbf{y} = \mathbf{0} \mathbf{y} = 0$$

contrary to  $\mathbf{x} \neq \mathbf{0}$ . Consequently there can be no  $\mathbf{y}$  such that  $A\mathbf{y} = \mathbf{x}$  and so  $A$  is not onto. This shows that 2.) implies 3.).

Finally, suppose 3.). If 1.) does not hold, then  $\det(A) \neq 0$  but then from Theorem 19.9.15 on Page 454  $A^{-1}$  exists and so for every  $\mathbf{y} \in \mathbb{R}^n$  there exists a unique  $\mathbf{x} \in \mathbb{R}^n$  such that  $A\mathbf{x} = \mathbf{y}$ . In fact  $\mathbf{x} = A^{-1}\mathbf{y}$ . Thus  $A$  would be onto contrary to 3.). This shows 3.) implies 1.) and proves the theorem.

**Corollary 19.9.23** *Let  $A$  be an  $n \times n$  matrix. Then the following are equivalent.*

1.  $\det(A) \neq 0$ .
2.  $A$  and  $A^T$  are one to one.
3.  $A$  is onto.

**Proof:** This follows immediately from the above theorem.

## 19.10 Exercises

1. Let  $m < n$  and let  $A$  be an  $m \times n$  matrix. Show that  $A$  is **not** one to one. **Hint:** Consider the  $n \times n$  matrix,  $A_1$  which is of the form

$$A_1 \equiv \begin{pmatrix} A \\ 0 \end{pmatrix}$$

where the 0 denotes an  $(n - m) \times n$  matrix of zeros. Thus  $\det A_1 = 0$  and so  $A_1$  is not one to one. Now observe that  $A_1 \mathbf{x}$  is the vector,

$$A_1 \mathbf{x} = \begin{pmatrix} A\mathbf{x} \\ \mathbf{0} \end{pmatrix}$$

which equals zero if and only if  $A\mathbf{x} = \mathbf{0}$ .

## 19.11 The Determinant And Volume

The determinant is the essential algebraic tool which provides a way to give a unified treatment of the concept of volume and it is this which is the most significant application of determinants.

**Lemma 19.11.1** Suppose  $A$  is an  $m \times n$  matrix where  $m > n$ . Then  $A$  does not map  $\mathbb{R}^n$  onto  $\mathbb{R}^m$ .

**Proof:** Suppose  $A$  did map  $\mathbb{R}^n$  onto  $\mathbb{R}^m$ . Then consider the  $m \times m$  matrix,

$$A_1 \equiv \begin{pmatrix} A & 0 \end{pmatrix}$$

where  $0$  refers to an  $n \times (m - n)$  matrix. Thus  $A_1$  cannot be onto  $\mathbb{R}^m$  because it has at least one column of zeros and so its determinant equals zero. However, if  $\mathbf{y} \in \mathbb{R}^m$  and  $A$  is onto, then there exists  $\mathbf{x} \in \mathbb{R}^n$  such that  $A\mathbf{x} = \mathbf{y}$ . Then

$$A_1 \begin{pmatrix} \mathbf{x} \\ \mathbf{0} \end{pmatrix} = A\mathbf{x} + \mathbf{0} = A\mathbf{x} = \mathbf{y}.$$

Since  $\mathbf{y}$  was arbitrary, it follows  $A_1$  would have to be onto.

The following proposition is a special case of a fundamental linear algebra result sometimes called the exchange theorem.

**Proposition 19.11.2** Suppose  $\{\mathbf{v}_1, \dots, \mathbf{v}_p\}$  are vectors in  $\mathbb{R}^n$  and  $\text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_p\} = \mathbb{R}^n$ . Then  $p \geq n$ .

**Proof:** Define a linear transformation from  $\mathbb{R}^p$  to  $\mathbb{R}^n$  as follows.

$$A\mathbf{x} \equiv \sum_{i=1}^p x_i \mathbf{v}_i.$$

(Why is this a linear transformation?) Thus  $A(\mathbb{R}^p) = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_p\} = \mathbb{R}^n$ . Then from the above lemma,  $p \geq n$  since if this is not so,  $A$  could not be onto.

**Proposition 19.11.3** Suppose  $\{\mathbf{v}_1, \dots, \mathbf{v}_p\}$  are vectors in  $\mathbb{R}^n$  such that  $p < n$ . Then there exist at least  $n - p$  vectors,  $\{\mathbf{w}_{p+1}, \dots, \mathbf{w}_n\}$  such that  $\mathbf{w}_i \cdot \mathbf{w}_j = \delta_{ij}$  and  $\mathbf{w}_k \cdot \mathbf{v}_j = 0$  for every  $j = 1, \dots, p$ .

**Proof:** Let  $A : \mathbb{R}^p \rightarrow \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_p\}$  be defined as in the above proposition so that  $A(\mathbb{R}^p) = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_p\}$ . Since  $p < n$  there exists  $\mathbf{z}_{p+1} \notin A(\mathbb{R}^p)$ . Then by Lemma 19.3.8 on Page 423 there exists  $\mathbf{x}_{p+1}$  such that

$$(\mathbf{z}_{p+1} - \mathbf{x}_{p+1}, A\mathbf{y}) = 0$$

for all  $\mathbf{y} \in \mathbb{R}^p$ . Let  $\mathbf{w}_{p+1} \equiv (\mathbf{z}_{p+1} - \mathbf{x}_{p+1}) / |\mathbf{z}_{p+1} - \mathbf{x}_{p+1}|$ . Now if  $p + 1 = n$ , stop.  $\{\mathbf{w}_{p+1}\}$  is the desired list of vectors. Otherwise, do for  $\text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_p, \mathbf{w}_{p+1}\}$  what was done for  $\text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_p\}$  using  $\mathbb{R}^{p+1}$  instead of  $\mathbb{R}^p$  and obtain  $\mathbf{w}_{p+2}$  in this way such that  $\mathbf{w}_{p+2} \cdot \mathbf{w}_{p+1} = 0$  and  $\mathbf{w}_{p+2} \cdot \mathbf{v}_k = 0$  for all  $k$ . Continue till a list of  $n - p$  vectors have been found.

Recall the geometric definition of the cross product of two vectors found on Page 322. As explained there, the magnitude of the cross product of two vectors was the area of the parallelogram determined by the two vectors. There was also a coordinate description of the cross product. In terms of the notation of Proposition 13.6.4 on Page 331 the  $i^{\text{th}}$  coordinate of the cross product is given by

$$\varepsilon_{ijk} u_j v_k$$

where the two vectors are  $(u_1, u_2, u_3)$  and  $(v_1, v_2, v_3)$ . Therefore, using the reduction identity of Lemma 13.6.3 on Page 331

$$\begin{aligned} |\mathbf{u} \times \mathbf{v}|^2 &= \varepsilon_{ijk} u_j v_k \varepsilon_{irs} u_r v_s \\ &= (\delta_{jr} \delta_{ks} - \delta_{kr} \delta_{js}) u_j v_k u_r v_s \\ &= u_j v_k u_j v_k - u_j v_k u_k v_j \\ &= (\mathbf{u} \cdot \mathbf{u})(\mathbf{v} \cdot \mathbf{v}) - (\mathbf{u} \cdot \mathbf{v})^2 \end{aligned}$$

which equals

$$\det \begin{pmatrix} \mathbf{u} \cdot \mathbf{u} & \mathbf{u} \cdot \mathbf{v} \\ \mathbf{u} \cdot \mathbf{v} & \mathbf{v} \cdot \mathbf{v} \end{pmatrix}.$$

Now recall the box product and how the box product was  $\pm$  the volume of the parallelepiped spanned by the three vectors. From the definition of the box product

$$\begin{aligned} \mathbf{u} \times \mathbf{v} \cdot \mathbf{w} &= \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ u_1 & u_2 & u_3 \\ v_1 & v_2 & v_3 \end{vmatrix} \cdot (w_1 \mathbf{i} + w_2 \mathbf{j} + w_3 \mathbf{k}) \\ &= \det \begin{pmatrix} w_1 & w_2 & w_3 \\ u_1 & u_2 & u_3 \\ v_1 & v_2 & v_3 \end{pmatrix}. \end{aligned}$$

Therefore,

$$|\mathbf{u} \times \mathbf{v} \cdot \mathbf{w}|^2 = \det \begin{pmatrix} w_1 & w_2 & w_3 \\ u_1 & u_2 & u_3 \\ v_1 & v_2 & v_3 \end{pmatrix}^2$$

which from the theory of determinants equals

$$\begin{aligned} &\det \begin{pmatrix} u_1 & u_2 & u_3 \\ v_1 & v_2 & v_3 \\ w_1 & w_2 & w_3 \end{pmatrix} \det \begin{pmatrix} u_1 & v_1 & w_1 \\ u_2 & v_2 & w_2 \\ u_3 & v_3 & w_3 \end{pmatrix} = \\ &\det \left( \begin{pmatrix} u_1 & u_2 & u_3 \\ v_1 & v_2 & v_3 \\ w_1 & w_2 & w_3 \end{pmatrix} \begin{pmatrix} u_1 & v_1 & w_1 \\ u_2 & v_2 & w_2 \\ u_3 & v_3 & w_3 \end{pmatrix} \right) = \\ &\det \begin{pmatrix} u_1^2 + u_2^2 + u_3^2 & u_1 v_1 + u_2 v_2 + u_3 v_3 & u_1 w_1 + u_2 w_2 + u_3 w_3 \\ u_1 v_1 + u_2 v_2 + u_3 v_3 & v_1^2 + v_2^2 + v_3^2 & v_1 w_1 + v_2 w_2 + v_3 w_3 \\ u_1 w_1 + u_2 w_2 + u_3 w_3 & v_1 w_1 + v_2 w_2 + v_3 w_3 & w_1^2 + w_2^2 + w_3^2 \end{pmatrix} \\ &= \det \begin{pmatrix} \mathbf{u} \cdot \mathbf{u} & \mathbf{u} \cdot \mathbf{v} & \mathbf{u} \cdot \mathbf{w} \\ \mathbf{u} \cdot \mathbf{v} & \mathbf{v} \cdot \mathbf{v} & \mathbf{v} \cdot \mathbf{w} \\ \mathbf{u} \cdot \mathbf{w} & \mathbf{v} \cdot \mathbf{w} & \mathbf{w} \cdot \mathbf{w} \end{pmatrix} \end{aligned}$$

You see there is a definite pattern emerging here. These earlier cases were for a parallelepiped determined by either two or three vectors in  $\mathbb{R}^3$ . It makes sense to speak of a parallelepiped in any number of dimensions.

**Definition 19.11.4** Let  $\mathbf{u}_1, \dots, \mathbf{u}_p$  be vectors in  $\mathbb{R}^k$ . The parallelepiped determined by these vectors will be denoted by  $P(\mathbf{u}_1, \dots, \mathbf{u}_p)$  and it is defined as

$$P(\mathbf{u}_1, \dots, \mathbf{u}_p) \equiv \left\{ \sum_{j=1}^p s_j \mathbf{u}_j : s_j \in [0, 1] \right\}.$$

The volume of this parallelepiped is defined as

$$\text{volume of } P(\mathbf{u}_1, \dots, \mathbf{u}_p) \equiv (\det(\mathbf{u}_i \cdot \mathbf{u}_j))^{1/2}.$$

In this definition,  $\mathbf{u}_i \cdot \mathbf{u}_j$  is the  $ij^{\text{th}}$  entry of a  $p \times p$  matrix. Note this definition agrees with all earlier notions of area and volume for parallelepipeds and it makes sense in any number of dimensions. However, it is important to verify the above determinant is nonnegative. After all, the above definition requires a square root of this determinant.

**Lemma 19.11.5** Let  $\mathbf{u}_1, \dots, \mathbf{u}_p$  be vectors in  $\mathbb{R}^k$  for some  $k$ . Then  $\det(\mathbf{u}_i \cdot \mathbf{u}_j) \geq 0$ .

**Proof:** Recall  $\mathbf{v} \cdot \mathbf{w} = \mathbf{v}^T \mathbf{w}$ . Therefore, in terms of matrix multiplication, the matrix  $(\mathbf{u}_i \cdot \mathbf{u}_j)$  is just the following

$$\overbrace{\begin{pmatrix} \mathbf{u}_1^T \\ \vdots \\ \mathbf{u}_p^T \end{pmatrix}}^{p \times k} \overbrace{\begin{pmatrix} \mathbf{u}_1 & \cdots & \mathbf{u}_p \end{pmatrix}}^{k \times p}$$

which is of the form

$$U^T U.$$

First show  $\det(U^T U) \geq 0$ . By Proposition 19.11.3 there are vectors,  $\mathbf{w}_{p+1}, \dots, \mathbf{w}_k$  such that  $\mathbf{w}_i \cdot \mathbf{w}_j = \delta_{ij}$  and for all  $i = 1, \dots, p$ , and  $l = p+1, \dots, k$ ,  $\mathbf{w}_l \cdot \mathbf{u}_i = 0$ . Then consider

$$U_1 = (\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{w}_{p+1}, \dots, \mathbf{w}_k) \equiv \begin{pmatrix} U & W \end{pmatrix}$$

where  $W^T W = I$ . Then

$$U_1^T U_1 = \begin{pmatrix} U^T \\ W^T \end{pmatrix} \begin{pmatrix} U & W \end{pmatrix} = \begin{pmatrix} U^T U & 0 \\ 0 & I \end{pmatrix}. \text{ (Why?)}$$

Now using the cofactor expansion method, this last  $k \times k$  matrix has determinant equal to  $\det(U^T U)$  (Why?) On the other hand this equals  $\det(U_1^T U_1) = \det(U_1) \det(U_1^T) = \det(U_1)^2 \geq 0$ .

In the case where  $k < p$ ,  $U^T U$  has the form  $W W^T$  where  $W = U^T$  has more rows than columns. Thus you can define the  $p \times p$  matrix,

$$W_1 \equiv \begin{pmatrix} W & 0 \end{pmatrix},$$

and in this case,

$$0 = \det W_1 W_1^T = \det \begin{pmatrix} W & 0 \end{pmatrix} \begin{pmatrix} W^T \\ 0 \end{pmatrix} = \det W W^T = \det U^T U.$$

This proves the lemma and shows the definition of volume is well defined.

Note it gives the right answer in the case where all the vectors are perpendicular. Here is why. Suppose  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  are vectors which have the property that  $\mathbf{u}_i \cdot \mathbf{u}_j = 0$  if  $i \neq j$ . Thus  $P(\mathbf{u}_1, \dots, \mathbf{u}_p)$  is a box which has all  $p$  sides perpendicular. What should its volume be? Shouldn't it equal the product of the lengths of the sides? What does  $\det(\mathbf{u}_i \cdot \mathbf{u}_j)$  give? The matrix  $(\mathbf{u}_i \cdot \mathbf{u}_j)$  is a diagonal matrix having the squares of the magnitudes of the sides down the diagonal. Therefore,  $\det(\mathbf{u}_i \cdot \mathbf{u}_j)^{1/2}$  equals the product of the lengths of the sides as it should. The matrix,  $(\mathbf{u}_i \cdot \mathbf{u}_j)$  whose determinant gives the square of the volume of the parallelepiped spanned by the vectors,  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  is called the Gramian matrix and sometimes the metric tensor.

These considerations are of great significance because they allow the computation in a systematic manner of  $k$  dimensional volumes of parallelepipeds which happen to be in  $\mathbb{R}^n$  for  $n \neq k$ . Think for example of a plane in  $\mathbb{R}^3$  and the problem of finding the area of something on this plane.

**Example 19.11.6** Find the equation of the plane containing the three points,  $(1, 2, 3)$ ,  $(0, 2, 1)$ , and  $(3, 1, 0)$ .

These three points determine two vectors, the one from  $(0, 2, 1)$  to  $(1, 2, 3)$ ,  $\mathbf{i} + 0\mathbf{j} + 2\mathbf{k}$ , and the one from  $(0, 2, 1)$  to  $(3, 1, 0)$ ,  $3\mathbf{i} + (-1)\mathbf{j} + (-1)\mathbf{k}$ . If  $(x, y, z)$  denotes a point in the plane, then the volume of the parallelepiped spanned by the vector from  $(0, 2, 1)$  to  $(x, y, z)$  and these other two vectors must be zero. Thus

$$\det \begin{pmatrix} x & y-2 & z-1 \\ 3 & -1 & -1 \\ 1 & 0 & 2 \end{pmatrix} = 0$$

Therefore,  $-2x - 7y + 13 + z = 0$  is the equation of the plane. You should check it contains all three points.

## 19.12 Exercises

1. Here are three vectors in  $\mathbb{R}^4$ :  $(1, 2, 0, 3)^T$ ,  $(2, 1, -3, 2)^T$ ,  $(0, 0, 1, 2)^T$ . Find the volume of the parallelepiped determined by these three vectors.
2. Here are two vectors in  $\mathbb{R}^4$ :  $(1, 2, 0, 3)^T$ ,  $(2, 1, -3, 2)^T$ . Find the volume of the parallelepiped determined by these two vectors.
3. Here are three vectors in  $\mathbb{R}^2$ :  $(1, 2)^T$ ,  $(2, 1)^T$ ,  $(0, 1)^T$ . Find the volume of the parallelepiped determined by these three vectors. Recall that from the above theorem, this should equal 0.
4. If there are  $n + 1$  or more vectors in  $\mathbb{R}^n$ , Lemma 19.11.5 implies the parallelepiped determined by these  $n + 1$  vectors must have zero volume. What is the geometric significance of this assertion?
5. Find the equation of the plane through the three points  $(1, 2, 3)$ ,  $(2, -3, 1)$ ,  $(1, 1, 7)$ .

## 19.13 Linear Systems Of Ordinary Differential Equations

Another important application of linear algebra is to linear systems of ordinary differential equations. Consider for  $t \in [a, b]$  the system

$$\mathbf{x}' = A(t)\mathbf{x}(t) + \mathbf{g}(t), \quad \mathbf{x}(c) = \mathbf{x}_0, \quad (19.54)$$

where  $c \in [a, b]$ ,  $A(t)$  is an  $n \times n$  matrix whose entries are continuous functions of  $t$ ,  $(a_{ij}(t))$  and  $\mathbf{g}(t)$  is a vector whose components are continuous functions of  $t$ . It turns out that this system satisfies the conditions of Theorem 15.7.6 on Page 370 with  $\mathbf{f}(t, \mathbf{x}) \equiv A(t)\mathbf{x} + \mathbf{g}(t)$ .

**Lemma 19.13.1** *The system, (19.54) satisfies the hypotheses of Theorem 15.7.6 on Page 370 with  $\mathbf{f}(t, \mathbf{x}) \equiv A(t)\mathbf{x} + \mathbf{g}(t)$ .*

**Proof:** First note that if the  $a_i$  are nonnegative numbers, then by the Cauchy Schwarz inequality,

$$\begin{aligned} \sum_{i=1}^n a_i &= \sum_{i=1}^n a_i 1 \leq \left( \sum_{i=1}^n a_i^2 \right)^{1/2} \left( \sum_{i=1}^n 1 \right)^{1/2} \\ &= n^{1/2} \left( \sum_{i=1}^n a_i^2 \right)^{1/2} \end{aligned}$$

so squaring both sides,

$$\left(\sum_{i=1}^n a_i\right)^2 \leq n \sum_{i=1}^n a_i^2.$$

Now let  $\mathbf{x} = (x_1, \dots, x_n)^T$  and  $\mathbf{x}_1 = (x_{11}, \dots, x_{1n})^T$ . Then letting

$$M \equiv \max \{|a_{ij}(t)| : t \in [a, b], i, j \leq n\},$$

it follows from the above inequality

$$\begin{aligned} |\mathbf{f}(t, \mathbf{x}) - \mathbf{f}(t, \mathbf{x}_1)| &= |A(t)(\mathbf{x} - \mathbf{x}_1)| \\ &= \left| \left( \sum_{i=1}^n \left| \sum_{j=1}^n a_{ij}(t)(x_j - x_{1j}) \right|^2 \right)^{1/2} \right| \\ &\leq M \left| \left( \sum_{i=1}^n \left( \sum_{j=1}^n |x_j - x_{1j}| \right)^2 \right)^{1/2} \right| \\ &\leq M \left| \left( \sum_{i=1}^n n \sum_{j=1}^n |x_j - x_{1j}|^2 \right)^{1/2} \right| \\ &= Mn \left( \sum_{j=1}^n |x_j - x_{1j}|^2 \right)^{1/2} = Mn |\mathbf{x} - \mathbf{x}_1|. \end{aligned}$$

Therefore, let  $K = Mn$ .

This lemma yields the following extremely important theorem.

**Theorem 19.13.2** *Let  $A(t)$  be a continuous  $n \times n$  matrix and let  $\mathbf{g}(t)$  be a continuous vector for  $t \in [a, b]$  and let  $c \in [a, b]$  and  $\mathbf{x}_0 \in \mathbb{R}^n$ . Then there exists a unique solution to (19.54) valid for  $t \in [a, b]$ .*

This theorem includes more examples of linear systems of equations than one typically encounters in any beginning course on differential equations. With this and the theory of determinants and matrices, some fundamental theorems about the nature of solutions to linear systems of equations become possible.

**Definition 19.13.3** *Let  $\Phi(t) = (\mathbf{x}_1(t), \dots, \mathbf{x}_n(t))$  be an  $n \times n$  matrix in which the  $i^{\text{th}}$  column is the vector,  $\mathbf{x}_i(t)$ . Then*

$$\Phi'(t) \equiv (\mathbf{x}'_1(t), \dots, \mathbf{x}'_n(t)).$$

$\Phi(t)$  is called a fundamental matrix for the equation

$$\mathbf{x}'(t) = A(t)\mathbf{x}(t), \quad (19.55)$$

on the interval  $[a, b]$  if for each  $i = 1, 2, \dots, n$ ,

$$\mathbf{x}'_i(t) = A(t)\mathbf{x}_i(t),$$

and for all  $t \in [a, b]$ ,

$$\Phi(t)^{-1} \text{ exists.}$$

A fundamental matrix consists of a matrix whose columns are solutions of (19.55). The following lemma is an amazing result about matrices of this form.

**Lemma 19.13.4** *Let  $\Phi(t)$  be a matrix which has the property that its columns are solutions of (19.55) where  $A(t)$  is a continuous  $n \times n$  matrix defined on an interval,  $[a, b]$ . Then if  $\mathbf{c}$  is any constant vector, it follows that  $\Phi(t)\mathbf{c}$  is a solution of (19.55). Furthermore, this gives all possible solutions to (19.55) if and only if  $\det(\Phi(t_0)) \neq 0$  for some  $t_0 \in [a, b]$ . (When all solutions to (19.55) are obtained in the form  $\Phi(t)\mathbf{c}$ ,  $\Phi(t)\mathbf{c}$  is called the general solution.) Finally,  $\det(\Phi(t))$  either equals zero for all  $t \in [a, b]$  or  $\det(\Phi(t))$  is never equal to zero for any  $t \in [a, b]$ .*

**Proof:** Let  $\Phi(t) = (\mathbf{x}_1(t), \dots, \mathbf{x}_n(t))$ . Then  $\Phi(t)\mathbf{c} = \sum_{i=1}^n c_i \mathbf{x}_i(t)$  and so

$$\begin{aligned} (\Phi(t)\mathbf{c})' &= \left( \sum_{i=1}^n c_i \mathbf{x}_i(t) \right)' = \sum_{i=1}^n c_i \mathbf{x}_i'(t) \\ &= \sum_{i=1}^n c_i A(t) \mathbf{x}_i(t) = A(t) \left( \sum_{i=1}^n c_i \mathbf{x}_i(t) \right) = A(t) (\Phi(t)\mathbf{c}). \end{aligned}$$

This proves the first part, also called the principle of superposition.

Suppose now that  $\det(\Phi(t_0)) \neq 0$  for some  $t_0 \in [a, b]$  and let  $\mathbf{z}$  be a solution of (19.55). Choose a constant vector,  $\mathbf{c}$  as follows.

$$\mathbf{c} \equiv \Phi(t_0)^{-1} \mathbf{z}(t_0).$$

Then  $\Phi(t_0)\mathbf{c} = \mathbf{z}(t_0)$  and defining  $\mathbf{x}(t) \equiv \Phi(t)\mathbf{c}$ , it follows  $\mathbf{x}$  is a solution to (19.55) which has the property that at the single point,  $t_0$ ,  $\mathbf{x}(t_0) = \mathbf{z}(t_0)$ . Therefore, by the uniqueness part of Theorem 19.13.2,  $\mathbf{x}(t) = \mathbf{z}(t)$  for all  $t \in [a, b]$ , showing that  $\mathbf{z}(t) = \Phi(t)\mathbf{c}$ . Since  $\mathbf{z}$  is an arbitrary solution to (19.55) this proves all solutions to the equation, (19.55), are obtained as  $\Phi(t)\mathbf{c}$  if  $\det(\Phi(t_0)) \neq 0$  at some  $t_0 \in [a, b]$ .

Now it is shown that if at some point,  $t_0$ ,  $\det(\Phi(t_0)) = 0$  then you do not obtain all solutions in this form. Since  $\det(\Phi(t_0)) = 0$ , it follows from Theorem 19.9.22 on Page 457 that  $\Phi(t_0)$  cannot map  $\mathbb{R}^n$  onto  $\mathbb{R}^n$  and so there exists  $\mathbf{z}_0 \in \mathbb{R}^n$  with the property that there is no solution,  $\mathbf{c}$ , to the equation

$$\Phi(t_0)\mathbf{c} = \mathbf{z}_0.$$

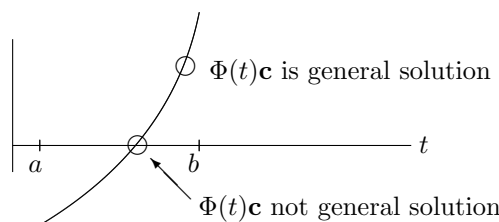
From the existence part of Theorem 19.13.2, there does exist a solution,  $\mathbf{z}$ , to the initial value problem

$$\mathbf{z}'(t) = A(t)\mathbf{z}(t), \quad \mathbf{z}(t_0) = \mathbf{z}_0,$$

and this solution cannot equal  $\Phi(t)\mathbf{c}$  for any choice of  $\mathbf{c}$  because if it did, you could plug in  $t_0$  and find that  $\Phi(t_0)\mathbf{c} = \mathbf{z}_0$ , a contradiction to how  $\mathbf{z}_0$  was chosen.

It remains to verify the last amazing assertion. This follows from the following picture which summarizes the above conclusions which have been proved at this point. The graph in the picture represents  $\det(\Phi(t))$  and shows it is impossible to have this function equal to zero at some point and yet nonzero at another.





Either  $\Phi(t)\mathbf{c}$  is the general solution or it is not. Therefore, the graph of  $\det(\Phi(t))$  cannot cross the  $t$  axis and this proves the lemma.

The function,  $t \rightarrow \det(\Phi(t))$  is known as the Wronskian<sup>2</sup>. and the above lemma is sometimes called the Wronskian alternative.

**Theorem 19.13.5** *Let  $A(t)$  be a continuous  $n \times n$  matrix and let  $c \in [a, b]$ . Then there exists a fundamental matrix for the equation,*

$$\mathbf{x}'(t) = A(t)\mathbf{x}(t),$$

$\Phi(t)$ , which satisfies  $\Phi(c) = I$ .

**Proof:** From Theorem 19.13.2 there exists  $\mathbf{x}_i$  satisfying (19.55) and the initial condition,  $\mathbf{x}_i(c) = \mathbf{e}_i$ . Now let  $\Phi(t) = (\mathbf{x}_1(t), \dots, \mathbf{x}_n(t))$ . Thus  $\Phi(c) = I$  and so  $\det(\Phi(t))$  is never equal to zero on  $[a, b]$  by the above lemma. Therefore, for all  $t \in [a, b]$ ,  $\Phi(t)^{-1}$  exists and consequently,  $\Phi(t)$  is a fundamental matrix. This proves the theorem.

The most important formula in the theory of differential equations is the variation of constants formula, sometimes called Green's formula. This formula is concerned with the solution to the system,

$$\mathbf{x}' = A(t)\mathbf{x}(t) + \mathbf{g}(t), \quad \mathbf{x}(c) = \mathbf{x}_0, \quad (19.56)$$

where  $c \in [a, b]$ ,  $A(t)$  is an  $n \times n$  matrix whose entries are continuous functions of  $t$ ,  $(a_{ij}(t))$  and  $\mathbf{g}(t)$  is a vector whose components are continuous functions of  $t$ . By Theorem 19.13.2 there exists a unique solution to this initial value problem but that is all this theorem says. The variation of constants formula gives a representation of the solution to this problem in terms of the fundamental matrix for  $\mathbf{x}' = A(t)\mathbf{x}$ . Look for a solution to (19.56) which is of the form  $\mathbf{x}(t) = \Phi(t)\mathbf{v}(t)$ ,  $\mathbf{v}(t) = (v_1(t), \dots, v_n(t))^T$ , where  $\Phi(t)$  is a fundamental matrix which satisfies  $\Phi(c) = I$ . Let  $\Phi(t) = (\mathbf{x}_1(t), \dots, \mathbf{x}_n(t))$ . Then

$$\Phi(t)\mathbf{v}(t) = \sum_{i=1}^n v_i(t)\mathbf{x}_i(t)$$

and so by the product rule,

$$\begin{aligned} (\Phi(t)\mathbf{v}(t))' &= \sum_{i=1}^n v_i'(t)\mathbf{x}_i(t) + v_i\mathbf{x}_i'(t) \\ &= \Phi(t)\mathbf{v}'(t) + \Phi'(t)\mathbf{v}(t). \end{aligned}$$

Therefore,  $\mathbf{x}$  will be a solution to the differential equation in (19.56) if and only if

$$\Phi(t)\mathbf{v}'(t) + \overbrace{\Phi'(t)\mathbf{v}(t)}^{=A(t)\Phi(t)\mathbf{v}(t)} = A(t)\Phi(t)\mathbf{v}(t) + \mathbf{g}(t).$$

<sup>2</sup>This is named after Wronski, a Polish mathematician who lived from 1778-1853.

Since  $\Phi(t)$  is a fundamental matrix,  $\Phi'(t) = A(t)\Phi(t)$  and so this equation reduces to

$$\Phi(t)\mathbf{v}'(t) = \mathbf{g}(t).$$

Now multiplying on both sides by  $\Phi(t)^{-1}$  which exists because  $\Phi(t)$  is a fundamental matrix, it follows that  $\mathbf{x}$  will be a solution to (19.56) if and only if

$$\mathbf{v}'(t) = \Phi(t)^{-1}\mathbf{g}(t).$$

From the formula for the inverse of a matrix,  $\Phi(t)^{-1}$  has all continuous entries and so the right side of the above is continuous. Therefore, it makes sense to define

$$\mathbf{v}(t) \equiv \int_c^t \Phi(s)^{-1}\mathbf{g}(s)ds + \mathbf{x}_0$$

and  $\Phi(t)\mathbf{v}(t)$  will be a solution to the differential equation of (19.56). Now consider

$$\mathbf{x}(t) = \Phi(t)\mathbf{v}(t) = \Phi(t)\mathbf{x}_0 + \Phi(t) \int_c^t \Phi(s)^{-1}\mathbf{g}(s)ds. \quad (19.57)$$

From the above discussion, it solves the differential equation of (19.56). However, it also satisfies the initial condition because  $\mathbf{x}(c) = \Phi(c)\mathbf{x}_0 + \Phi(c)\mathbf{0} = I\mathbf{x}_0$ . Therefore, by Theorem 19.13.2 the formula given in (19.57) is the solution to the initial value problem (19.56). It is (19.57) which is called the variation of constants formula.

## 19.14 Exercises

1. The differential equation of undamped oscillation is  $y'' + \omega^2 y = 0$  where  $\omega \neq 0$ . Define a new variable,  $z = y'$  and show this differential equation may be written as the first order system

$$\begin{pmatrix} z \\ y \end{pmatrix}' = \begin{pmatrix} 0 & -\omega^2 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} z \\ y \end{pmatrix}.$$

Now show a fundamental matrix for this equation is

$$\begin{pmatrix} \omega \cos \omega t & -\omega \sin \omega t \\ \sin \omega t & \cos \omega t \end{pmatrix}.$$

Finally, show that all solutions of the equation,  $y'' + \omega^2 y = 0$  are of the form  $c_1 \cos \omega t + c_2 \sin \omega t$ .

2. Do the same thing for the equation  $y'' - \omega^2 y = 0$  that was done in Problem 1.
3. Now suppose you have the equation,  $y'' + 2by' + 4cy = 0$ . Try to reduce this to one of the above problems by changing the equation to one for the dependent variable  $z = e^{\lambda t}y$  and choosing  $\lambda$  appropriately.
4. Show that any differential equation of the form  $y^{(n)} + a_1(t)y^{(n-1)} + \cdots + a_{n-1}(t)y' + a_n(t)y = 0$  can be written as a first order system using a technique similar to that outlined in Problem 1.
5. Show that if  $\Phi(t)$  is a fundamental matrix for  $\mathbf{x}' = A(t)\mathbf{x}$  and  $\mathbf{x}'_p = A(t)\mathbf{x}_p + \mathbf{f}(t)$  then if  $\mathbf{z}' = A(t)\mathbf{z} + \mathbf{f}(t)$ , there exists  $\mathbf{c}$  a constant vector such that  $\mathbf{x}_p + \Phi(t)\mathbf{c} = \mathbf{z}$ . Thus the general solution for the equation  $\mathbf{x}' = A(t)\mathbf{x} + \mathbf{f}(t)$  is  $\mathbf{x}_p + \Phi(t)\mathbf{c}$  where  $\mathbf{c}$  is an arbitrary constant vector. The solution,  $\mathbf{x}_p$  is called a particular solution.

6. He has forgotten the continuous matrix,  $A(t)$  for  $t \in [-1, 1]$  used in the differential equation  $\mathbf{x}' = A(t)\mathbf{x}$  but remembers there were two solutions.

$$\mathbf{x}_1(t) = \begin{pmatrix} \sin t \\ t^2 \end{pmatrix} \text{ and } \mathbf{x}_2(t) = \begin{pmatrix} te^t \\ t \end{pmatrix}$$

Find the continuous matrix if possible. If not possible, tell why.

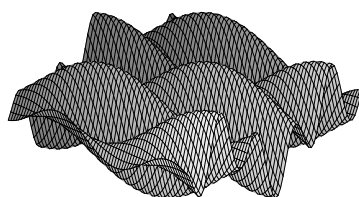


# Functions Of Many Variables

## 20.1 The Graph Of A Function Of Two Variables

With vector valued functions of many variables, it doesn't take long before it is impossible to draw meaningful pictures. This is because one needs more than three dimensions to accomplish the task and we can only visualize things in three dimensions. Ultimately, one of the main purposes of calculus is to free us from the tyranny of art. In calculus, we are permitted and even required to think in a meaningful way about things which cannot be drawn. However, it is certainly interesting to consider some things which can be visualized and this will help to formulate and understand more general notions which make sense in contexts which cannot be visualized. One of these is the concept of a scalar valued function of two variables.

Let  $f(x, y)$  denote a scalar valued function of two variables evaluated at the point  $(x, y)$ . Its graph consists of the set of points,  $(x, y, z)$  such that  $z = f(x, y)$ . How does one go about depicting such a graph? The usual way is to fix one of the variables, say  $x$  and consider the function  $z = f(x, y)$  where  $y$  is allowed to vary and  $x$  is fixed. Graphing this would give a curve which lies in the surface to be depicted. Then do the same thing for other values of  $x$  and the result would depict the graph desired graph. Computers do this very well. The following is the graph of the function  $z = \cos(x) \sin(2x + y)$  drawn using Maple, a computer algebra system.<sup>1</sup>

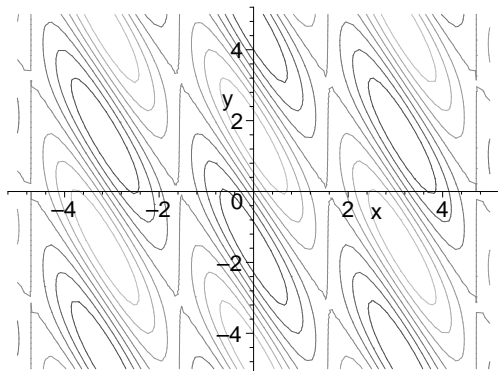


Notice how elaborate this picture is. The lines in the drawing correspond to taking one of the variables constant and graphing the curve which results. The computer did this drawing in seconds but you couldn't do it as well if you spent all day on it. I used a grid consisting of 70 choices for  $x$  and 70 choices for  $y$ .

Sometimes attempts are made to understand three dimensional objects like the above graph by looking at contour graphs in two dimensions. The contour graph of the above three dimensional graph is below and comes from using the computer algebra system again.

---

<sup>1</sup>I used Maple and exported the graph as an eps. file which I then imported into this document.

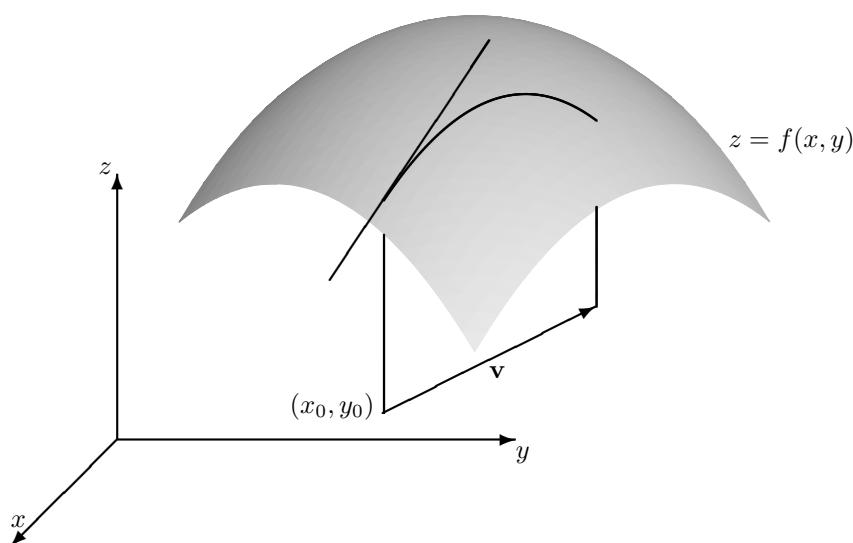


This is in two dimensions and the different lines in two dimensions correspond to points on the three dimensional graph which have the same  $z$  value. If you have looked at a weather map, these lines are called isotherms or isobars depending on whether the function involved is temperature or pressure. In a contour geographic map, the contour lines represent constant altitude. If many contour lines are close to each other, this indicates rapid change in the altitude, temperature, pressure, or whatever else may be measured.

A scalar function of three variables, cannot be visualized because four dimensions are required. However, some people like to try and visualize even these examples. This is done by looking at level surfaces in  $\mathbb{R}^3$  which are defined as surfaces where the function assumes a constant value. They play the role of contour lines for a function of two variables. As a simple example, consider  $f(x, y, z) = x^2 + y^2 + z^2$ . The level surfaces of this function would be concentric spheres centered at  $\mathbf{0}$ . (Why?) Another way to visualize objects in higher dimensions involves the use of color and animation. However, there really are limits to what you can accomplish in this direction. So much for art.

## 20.2 The Directional Derivative

The directional derivative is just what its name suggests. It is the derivative of a function in a particular direction. The following picture illustrates the situation in the case of a function of two variables.



In this picture,  $\mathbf{v} \equiv (v_1, v_2)$  is a unit vector in the  $xy$  plane and  $\mathbf{x}_0 \equiv (x_0, y_0)$  is a point in the  $xy$  plane. When  $(x, y)$  moves in the direction of  $\mathbf{v}$ , this results in a change in  $z = f(x, y)$  as shown in the picture. The directional derivative in this direction is defined as

$$\lim_{t \rightarrow 0} \frac{f(x_0 + tv_1, y_0 + tv_2) - f(x_0, y_0)}{t}.$$

It tells how fast  $z$  is changing in this direction. If you looked at it from the side, you would be getting the slope of the indicated tangent line. A simple example of this is a person climbing a mountain. He could go various directions, some steeper than others. The directional derivative is just a measure of the steepness in a given direction. This motivates the following general definition of the directional derivative.

**Definition 20.2.1** Let  $f : U \rightarrow \mathbb{R}$  where  $U$  is an open set in  $\mathbb{R}^n$  and let  $\mathbf{v}$  be a unit vector. For  $\mathbf{x} \in U$ , define the directional derivative of  $f$  in the direction,  $\mathbf{v}$ , at the point  $\mathbf{x}$  as

$$D_{\mathbf{v}}f(\mathbf{x}) \equiv \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t}.$$

**Example 20.2.2** Find the directional derivative of the function,  $f(x, y) = x^2y$  in the direction of  $\mathbf{i} + \mathbf{j}$  at the point  $(1, 2)$ .

First you need a unit vector which has the same direction as the given vector. This unit vector is  $\mathbf{v} \equiv \left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}\right)$ . Then to find the directional derivative from the definition, write the difference quotient described above. Thus  $f(\mathbf{x} + t\mathbf{v}) = \left(1 + \frac{t}{\sqrt{2}}\right)^2 \left(2 + \frac{t}{\sqrt{2}}\right)$  and  $f(\mathbf{x}) = 2$ . Therefore,

$$\frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} = \frac{\left(1 + \frac{t}{\sqrt{2}}\right)^2 \left(2 + \frac{t}{\sqrt{2}}\right) - 2}{t},$$

and to find the directional derivative, you take the limit of this as  $t \rightarrow 0$ . However, this difference quotient equals  $\frac{1}{4}\sqrt{2}(10 + 4t\sqrt{2} + t^2)$  and so, letting  $t \rightarrow 0$ ,

$$D_{\mathbf{v}}f(1, 2) = \left(\frac{5}{2}\sqrt{2}\right).$$

There is something you must keep in mind about this. The direction vector must always be a unit vector<sup>2</sup>.

There are some special unit vectors which come to mind immediately. These are the vectors,  $\mathbf{e}_i$  where

$$\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)^T$$

and the 1 is in the  $i^{\text{th}}$  position.

**Definition 20.2.3** Let  $U$  be an open subset of  $\mathbb{R}^n$  and let  $f : U \rightarrow \mathbb{R}$ . Then letting  $\mathbf{x} = (x_1, \dots, x_n)^T$  be a typical element of  $\mathbb{R}^n$ ,

$$\frac{\partial f}{\partial x_i}(\mathbf{x}) \equiv D_{\mathbf{e}_i}f(\mathbf{x}).$$

This is called the partial derivative of  $f$ . Thus,

$$\begin{aligned} \frac{\partial f}{\partial x_i}(\mathbf{x}) &\equiv \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{e}_i) - f(\mathbf{x})}{t} \\ &= \lim_{t \rightarrow 0} \frac{f(x_1, \dots, x_i + t, \dots, x_n) - f(x_1, \dots, x_i, \dots, x_n)}{t}, \end{aligned}$$

<sup>2</sup>Actually, there is a more general formulation of this known as the Gateaux derivative in which the length of  $\mathbf{v}$  is not considered but it will not be considered here.

and to find the partial derivative, differentiate with respect to the variable of interest and regard all the others as constants. Other notation for this partial derivative is  $f_{x_i}$ ,  $f_{,i}$ , or  $D_i f$ . If  $y = f(\mathbf{x})$ , the partial derivative of  $f$  with respect to  $x_i$  may also be denoted by

$$\frac{\partial y}{\partial x_i} \text{ or } y_{x_i}.$$

**Example 20.2.4** Find  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ , and  $\frac{\partial f}{\partial z}$  if  $f(x, y) = y \sin x + x^2 y + z$ .

From the definition above,  $\frac{\partial f}{\partial x} = y \cos x + 2xy$ ,  $\frac{\partial f}{\partial y} = \sin x + x^2$ , and  $\frac{\partial f}{\partial z} = 1$ . Having taken one partial derivative, there is no reason to stop doing it. Thus, one could take the partial derivative with respect to  $y$  of the partial derivative with respect to  $x$ , denoted by  $\frac{\partial^2 f}{\partial y \partial x}$  or  $f_{xy}$ . In the above example,

$$\frac{\partial^2 f}{\partial y \partial x} = f_{xy} = \cos x + 2x.$$

Also observe that

$$\frac{\partial^2 f}{\partial x \partial y} = f_{yx} = \cos x + 2x.$$

Higher order partial derivatives are defined by analogy to the above. Thus in the above example,

$$f_{yxx} = -\sin x + 2.$$

These partial derivatives,  $f_{xy}$  are called mixed partial derivatives.

The concept of a directional derivative for a vector valued function is also easy to define although the geometric significance expressed in pictures is not.

**Definition 20.2.5** Let  $\mathbf{f} : U \rightarrow \mathbb{R}^p$  where  $U$  is an open set in  $\mathbb{R}^n$  and let  $\mathbf{v}$  be a unit vector. For  $\mathbf{x} \in U$ , define the directional derivative of  $\mathbf{f}$  in the direction,  $\mathbf{v}$ , at the point  $\mathbf{x}$  as

$$D_{\mathbf{v}} \mathbf{f}(\mathbf{x}) \equiv \lim_{t \rightarrow 0} \frac{\mathbf{f}(\mathbf{x} + t\mathbf{v}) - \mathbf{f}(\mathbf{x})}{t}.$$

**Example 20.2.6** Let  $\mathbf{f}(x, y) = (xy^2, yx)^T$ . Find the directional derivative in the direction  $(1, 2)^T$  at the point  $(x, y)$ .

First, a unit vector in this direction is  $(1/\sqrt{5}, 2/\sqrt{5})^T$  and from the definition, the desired limit is

$$\begin{aligned} & \lim_{t \rightarrow 0} \frac{\left( (x + t(1/\sqrt{5})) (y + t(2/\sqrt{5}))^2 - xy^2, (x + t(1/\sqrt{5})) (y + t(2/\sqrt{5})) - xy \right)}{t} \\ &= \lim_{t \rightarrow 0} \left( \frac{4}{5}xy\sqrt{5} + \frac{4}{5}xt + \frac{1}{5}\sqrt{5}y^2 + \frac{4}{5}ty + \frac{4}{25}t^2\sqrt{5}, \frac{2}{5}x\sqrt{5} + \frac{1}{5}y\sqrt{5} + \frac{2}{5}t \right) \\ &= \left( \frac{4}{5}xy\sqrt{5} + \frac{1}{5}\sqrt{5}y^2, \frac{2}{5}x\sqrt{5} + \frac{1}{5}y\sqrt{5} \right). \end{aligned}$$

You see from this example and the above definition that all you have to do is to form the vector which is obtained by replacing each component of the vector with its directional derivative. In particular, you can take partial derivatives of vector valued functions and use the same notation.

**Example 20.2.7** Find the partial derivative with respect to  $x$  of the function  $\mathbf{f}(x, y, z, w) = (xy^2, z \sin(xy), z^3x)^T$ .

From the above definition,  $\mathbf{f}_x(x, y, z) = D_1 \mathbf{f}(x, y, z) = (y^2, zy \cos(xy), z^3)^T$ .



## 20.3 Exercises

1. Why must the vector in the definition of the directional derivative be a unit vector?  
**Hint:** Suppose not. Would the directional derivative be a correct manifestation of steepness?

2. Find  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ , and  $\frac{\partial f}{\partial z}$  for  $f =$

- (a)  $x^2y + \cos(xy) + z^3y$
- (b)  $e^{x^2+y^2}z \sin(x+y)$
- (c)  $z^2 \sin^3(e^{x^2+y^3})$
- (d)  $x^2 \cos(\sin(\tan(z^2 + y^2)))$
- (e)  $x^{y^2+z}$

3. Suppose

$$f(x, y) = \begin{cases} \frac{2xy + 6x^3 + 12xy^2 + 18y^3 + \sin(x^3) + \tan(3y^3)}{3x^2 + 6y^2} & \text{if } (x, y) \neq (0, 0) \\ 0 & \text{if } (x, y) = (0, 0). \end{cases}$$

Find  $\frac{\partial f}{\partial x}(0, 0)$  and  $\frac{\partial f}{\partial y}(0, 0)$ .

4. Find  $f_x, f_y, f_z, f_{xy}, f_{yx}, f_{xz}, f_{zx}, f_{zy}, f_{yz}$  for the following and form a conjecture about the mixed partial derivatives.

- (a)  $x^2y^3z^4 + \sin(xyz)$
- (b)  $\sin(xyz) + x^2yz$
- (c)  $z \ln|x^2 + y^2 + 1|$
- (d)  $e^{x^2+y^2+z^2}$
- (e)  $\tan(xyz)$

5. Suppose  $f : U \rightarrow \mathbb{R}$  where  $U$  is an open set and suppose that  $\mathbf{x} \in U$  has the property that for all  $\mathbf{y}$  near  $\mathbf{x}$ ,  $f(\mathbf{x}) \leq f(\mathbf{y})$ . Prove that if  $f$  has all of its partial derivatives at  $\mathbf{x}$ , then  $f_{x_i}(\mathbf{x}) = 0$  for each  $x_i$ . **Hint:** This is just a repeat of the similar one variable theorem given earlier, Theorem 6.5.2 on Page 126. You just do this argument given earlier for each variable to get the conclusion.

6. As an important application of Problem 5 consider the following. Experiments are done at  $n$  times,  $t_1, t_2, \dots, t_n$  and at each time there results a collection of numerical outcomes. Denote by  $\{(t_i, x_i)\}_{i=1}^p$  the set of all such pairs and try to find numbers  $a$  and  $b$  such that the line  $x = at + b$  approximates these ordered pairs as well as possible in the sense that out of all choices of  $a$  and  $b$ ,  $\sum_{i=1}^p (at_i + b - x_i)^2$  is as small as possible. In other words, you want to minimize the function of two variables,  $f(a, b) \equiv \sum_{i=1}^p (at_i + b - x_i)^2$ . Find a formula for  $a$  and  $b$  in terms of the given ordered pairs. You will be finding the formula for the least squares regression line.

## 20.4 Mixed Partial Derivatives

Under certain conditions the mixed partial derivatives will always be equal. This astonishing fact is due to Euler in 1734.

**Theorem 20.4.1** Suppose  $f : U \subseteq \mathbb{R}^2 \rightarrow \mathbb{R}$  where  $U$  is an open set on which  $f_x, f_y, f_{xy}$  and  $f_{yx}$  exist. Then if  $f_{xy}$  and  $f_{yx}$  are continuous at the point  $(x, y) \in U$ , it follows

$$f_{xy}(x, y) = f_{yx}(x, y).$$

**Proof:** Since  $U$  is open, there exists  $r > 0$  such that  $B((x, y), r) \subseteq U$ . Now let  $|t|, |s| < r/2$  and consider

$$\Delta(s, t) \equiv \frac{1}{st} \left\{ \overbrace{f(x+t, y+s) - f(x+t, y)}^{h(t)} - \overbrace{(f(x, y+s) - f(x, y))}^{h(0)} \right\}. \quad (20.1)$$

Note that  $(x+t, y+s) \in U$  because

$$\begin{aligned} |(x+t, y+s) - (x, y)| &= |(t, s)| = (t^2 + s^2)^{1/2} \\ &\leq \left( \frac{r^2}{4} + \frac{r^2}{4} \right)^{1/2} = \frac{r}{\sqrt{2}} < r. \end{aligned}$$

As implied above,  $h(t) \equiv f(x+t, y+s) - f(x+t, y)$ . Therefore, by the mean value theorem from calculus and the (one variable) chain rule,

$$\begin{aligned} \Delta(s, t) &= \frac{1}{st} (h(t) - h(0)) = \frac{1}{st} h'(\alpha t) t \\ &= \frac{1}{s} (f_x(x + \alpha t, y+s) - f_x(x + \alpha t, y)) \end{aligned}$$

for some  $\alpha \in (0, 1)$ . Applying the mean value theorem again,

$$\Delta(s, t) = f_{xy}(x + \alpha t, y + \beta s)$$

where  $\alpha, \beta \in (0, 1)$ .

If the terms  $f(x+t, y)$  and  $f(x, y+s)$  are interchanged in (20.1),  $\Delta(s, t)$  is also unchanged and the above argument shows there exist  $\gamma, \delta \in (0, 1)$  such that

$$\Delta(s, t) = f_{yx}(x + \gamma t, y + \delta s).$$

Letting  $(s, t) \rightarrow (0, 0)$  and using the continuity of  $f_{xy}$  and  $f_{yx}$  at  $(x, y)$ ,

$$\lim_{(s,t) \rightarrow (0,0)} \Delta(s, t) = f_{xy}(x, y) = f_{yx}(x, y).$$

This proves the theorem.

The following is obtained from the above by simply fixing all the variables except for the two of interest.

**Corollary 20.4.2** Suppose  $U$  is an open subset of  $\mathbb{R}^n$  and  $f : U \rightarrow \mathbb{R}$  has the property that for two indices,  $k, l$ ,  $f_{x_k}, f_{x_l}, f_{x_l x_k}$ , and  $f_{x_k x_l}$  exist on  $U$  and  $f_{x_k x_l}$  and  $f_{x_l x_k}$  are both continuous at  $\mathbf{x} \in U$ . Then  $f_{x_k x_l}(\mathbf{x}) = f_{x_l x_k}(\mathbf{x})$ .

It is necessary to assume the mixed partial derivatives are continuous in order to assert they are equal. The following is a well known example [1].

**Example 20.4.3** *Let*

$$f(x, y) = \begin{cases} \frac{xy(x^2 - y^2)}{x^2 + y^2} & \text{if } (x, y) \neq (0, 0) \\ 0 & \text{if } (x, y) = (0, 0) \end{cases}$$

From the definition of partial derivatives it follows immediately that  $f_x(0, 0) = f_y(0, 0) = 0$ . Using the standard rules of differentiation, for  $(x, y) \neq (0, 0)$ ,

$$f_x = y \frac{x^4 - y^4 + 4x^2y^2}{(x^2 + y^2)^2}, \quad f_y = x \frac{x^4 - y^4 - 4x^2y^2}{(x^2 + y^2)^2}$$

Now

$$\begin{aligned} f_{xy}(0, 0) &\equiv \lim_{y \rightarrow 0} \frac{f_x(0, y) - f_x(0, 0)}{y} \\ &= \lim_{y \rightarrow 0} \frac{-y^4}{(y^2)^2} = -1 \end{aligned}$$

while

$$\begin{aligned} f_{yx}(0, 0) &\equiv \lim_{x \rightarrow 0} \frac{f_y(x, 0) - f_y(0, 0)}{x} \\ &= \lim_{x \rightarrow 0} \frac{x^4}{(x^2)^2} = 1 \end{aligned}$$

showing that although the mixed partial derivatives do exist at  $(0, 0)$ , they are not equal there.

## 20.5 The Limit Of A Function Of Many Variables

Limits of scalar valued functions of one variable and later vector valued functions of one variable were considered earlier. Now consider vector valued functions of many variables. When  $\mathbf{f} : D(\mathbf{f}) \rightarrow \mathbb{R}^q$  one can only consider in a meaningful way limits at limit points of the set,  $D(\mathbf{f})$  and this concept is defined next.

**Definition 20.5.1** *Let  $A$  denote a nonempty subset of  $\mathbb{R}^p$ . A point,  $\mathbf{x}$  is said to be a limit point of the set,  $A$  if for every  $r > 0$ ,  $B(\mathbf{x}, r)$  contains infinitely many points of  $A$ .*

**Lemma 20.5.2** *If  $A$  is a nonempty open set in  $\mathbb{R}^p$ , then every point of  $A$  is a limit point of  $A$ .*

**Proof:** Let  $\mathbf{x} \in A$  and let  $B(\mathbf{x}, r)$  be a ball as above. Since  $\mathbf{x} \in A$ , an open set,  $B(\mathbf{x}, r_1) \subseteq A$  for some  $r_1 > 0$ . Then letting  $\delta = \min(r_1, r)$ , it follows  $B(\mathbf{x}, \delta) \subseteq A \cap B(\mathbf{x}, r)$ . Now every nonempty ball in  $\mathbb{R}^p$  contains infinitely many points (Why?) and this proves the lemma.

The case where  $A$  is open will be the one of most interest but many other sets have limit points.

**Definition 20.5.3** *Let  $\mathbf{f} : D(\mathbf{f}) \subseteq \mathbb{R}^p \rightarrow \mathbb{R}^q$  where  $q, p \geq 1$  be a function and let  $\mathbf{x}$  be a limit point of  $D(\mathbf{f})$ . Then*

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L}$$

if and only if the following condition holds. For all  $\varepsilon > 0$  there exists  $\delta > 0$  such that if

$$0 < |\mathbf{y} - \mathbf{x}| < \delta \text{ and } \mathbf{y} \in D(\mathbf{f})$$

then,

$$|\mathbf{L} - \mathbf{f}(\mathbf{y})| < \varepsilon.$$

**Proposition 20.5.4** Let  $\mathbf{f} : D(\mathbf{f}) \subseteq \mathbb{R}^p \rightarrow \mathbb{R}^q$  where  $q, p \geq 1$  be a function and let  $\mathbf{x}$  be a limit point of  $D(\mathbf{f})$ . Then if  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y})$  exists, it must be unique.

**Proof:** Suppose  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L}_1$  and  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L}_2$ . Then for  $\varepsilon > 0$  be given, let  $\delta_i > 0$  correspond to  $\mathbf{L}_i$  in the definition of the limit and let  $\delta = \min(\delta_1, \delta_2)$ . Since  $\mathbf{x}$  is a limit point, there exists  $\mathbf{y} \in B(\mathbf{x}, \delta) \cap D(\mathbf{f})$ . Therefore,

$$\begin{aligned} |\mathbf{L}_1 - \mathbf{L}_2| &\leq |\mathbf{L}_1 - \mathbf{f}(\mathbf{y})| + |\mathbf{f}(\mathbf{y}) - \mathbf{L}_2| \\ &< \varepsilon + \varepsilon = 2\varepsilon. \end{aligned}$$

Since  $\varepsilon > 0$  is arbitrary, this shows  $\mathbf{L}_1 = \mathbf{L}_2$ .

**Theorem 20.5.5** Suppose  $\mathbf{x}$  is a limit point of  $D(\mathbf{f})$  and  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L}$ ,  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{g}(\mathbf{y}) = \mathbf{K}$  where  $\mathbf{K}$  and  $\mathbf{L}$  are vectors in  $\mathbb{R}^p$  for  $p \geq 1$ . Then if  $a, b \in \mathbb{R}$ ,

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} a\mathbf{f}(\mathbf{y}) + b\mathbf{g}(\mathbf{y}) = a\mathbf{L} + b\mathbf{K}, \quad (20.2)$$

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f} \cdot \mathbf{g}(\mathbf{y}) = \mathbf{L} \cdot \mathbf{K} \quad (20.3)$$

Also, if  $\mathbf{h}$  is a continuous function defined near  $\mathbf{L}$ , then

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{h} \circ \mathbf{f}(\mathbf{y}) = \mathbf{h}(\mathbf{L}). \quad (20.4)$$

For a vector valued function,  $\mathbf{f}(\mathbf{y}) = (f_1(\mathbf{y}), \dots, f_q(\mathbf{y}))$ ,  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L} = (L_1, \dots, L_k)^T$  if and only if

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} f_k(\mathbf{y}) = L_k \quad (20.5)$$

for each  $k = 1, \dots, p$ .

**Proof:** The proof of (20.2) is left for you. It is like a corresponding theorem for continuous functions. Consider (20.3). Let  $\varepsilon > 0$  be given. Then by the triangle inequality, the properties of the dot product and the Cauchy Schwartz inequality,

$$\begin{aligned} |\mathbf{f} \cdot \mathbf{g}(\mathbf{y}) - \mathbf{L} \cdot \mathbf{K}| &\leq |\mathbf{f} \cdot \mathbf{g}(\mathbf{y}) - \mathbf{f}(\mathbf{y}) \cdot \mathbf{K}| + |\mathbf{f}(\mathbf{y}) \cdot \mathbf{K} - \mathbf{L} \cdot \mathbf{K}| \\ &\leq |\mathbf{f}(\mathbf{y})| |\mathbf{g}(\mathbf{y}) - \mathbf{K}| + |\mathbf{K}| |\mathbf{f}(\mathbf{y}) - \mathbf{L}|. \end{aligned} \quad (20.6)$$

There exists  $\delta_1$  such that if  $0 < |\mathbf{y} - \mathbf{x}| < \delta_1$ , then

$$|\mathbf{f}(\mathbf{y}) - \mathbf{L}| < 1,$$

and so for such  $\mathbf{y}$ , it follows from the triangle inequality that  $|\mathbf{f}(\mathbf{y})| < 1 + |\mathbf{L}|$ . Therefore, for  $0 < |\mathbf{y} - \mathbf{x}| < \delta_1$ ,

$$|\mathbf{f} \cdot \mathbf{g}(\mathbf{y}) - \mathbf{L} \cdot \mathbf{K}| \leq (1 + |\mathbf{K}| + |\mathbf{L}|) [|\mathbf{g}(\mathbf{y}) - \mathbf{K}| + |\mathbf{f}(\mathbf{y}) - \mathbf{L}|]. \quad (20.7)$$

Now let  $0 < \delta_2$  be such that for  $0 < |\mathbf{x} - \mathbf{y}| < \delta_2$ ,

$$|\mathbf{f}(\mathbf{y}) - \mathbf{L}| < \frac{\varepsilon}{2(1 + |\mathbf{K}| + |\mathbf{L}|)}, \quad |\mathbf{g}(\mathbf{y}) - \mathbf{K}| < \frac{\varepsilon}{2(1 + |\mathbf{K}| + |\mathbf{L}|)}.$$

Then letting  $0 < \delta \leq \min(\delta_1, \delta_2)$ , it follows from (20.7) that

$$|\mathbf{f} \cdot \mathbf{g}(\mathbf{y}) - \mathbf{L} \cdot \mathbf{K}| < \varepsilon$$

and this proves (20.3).

Consider (20.4). Since  $\mathbf{h}$  is continuous near  $\mathbf{L}$ , it follows that for  $\varepsilon > 0$  given, there exists  $\eta > 0$  such that if  $|\mathbf{y} - \mathbf{L}| < \eta$ , then

$$|\mathbf{h}(\mathbf{y}) - \mathbf{h}(\mathbf{L})| < \varepsilon$$

Since  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L}$ , there exists  $\delta > 0$  such that if  $0 < |\mathbf{y} - \mathbf{x}| < \delta$ , then

$$|\mathbf{f}(\mathbf{y}) - \mathbf{L}| < \eta.$$

Therefore, if  $0 < |\mathbf{y} - \mathbf{x}| < \delta$ , then

$$|\mathbf{h}(\mathbf{f}(\mathbf{y})) - \mathbf{h}(\mathbf{L})| < \varepsilon.$$

Consider (20.5). Suppose first that  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L}$ . Then if  $\varepsilon > 0$  is given, there exists  $\delta > 0$  such that if  $|\mathbf{y} - \mathbf{x}| < \delta$ ,

$$|\mathbf{f}(\mathbf{y}) - \mathbf{L}| < \varepsilon.$$

But for each  $k$ ,

$$|f_k(\mathbf{y}) - L_k| \leq |\mathbf{f}(\mathbf{y}) - \mathbf{L}|$$

and so if  $|\mathbf{y} - \mathbf{x}| < \delta$ ,

$$|f_k(\mathbf{y}) - L_k| \leq |\mathbf{f}(\mathbf{y}) - \mathbf{L}| < \varepsilon.$$

Next suppose  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} f_k(\mathbf{y}) = L_k$  for each  $k$ . Then if  $\varepsilon > 0$  is given, there exists  $\delta_k$  such that if  $0 < |\mathbf{y} - \mathbf{x}| < \delta_k$ ,

$$|f_k(\mathbf{y}) - L_k| < \frac{\varepsilon}{\sqrt{q}}.$$

Then, letting  $0 < \delta \leq \min(\delta_1, \dots, \delta_q)$ , it follows that if  $0 < |\mathbf{y} - \mathbf{x}| < \delta$ ,

$$\begin{aligned} |\mathbf{f}(\mathbf{y}) - \mathbf{L}| &\equiv \left( \sum_{k=1}^q |f_k(\mathbf{y}) - L_k|^2 \right)^{1/2} \\ &< \left( \sum_{k=1}^q \frac{\varepsilon^2}{q} \right)^{1/2} = \varepsilon. \end{aligned}$$

This proves the theorem.

**Theorem 20.5.6** For  $\mathbf{f} : D(\mathbf{f}) \rightarrow \mathbb{R}^q$  and  $\mathbf{x} \in D(\mathbf{f})$  such that  $\mathbf{x}$  is a limit point of  $D(\mathbf{f})$ , it follows  $\mathbf{f}$  is continuous at  $\mathbf{x}$  if and only if  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{f}(\mathbf{x})$ .

**Proof:** It follows immediately that if  $\mathbf{f}$  is continuous at  $\mathbf{x}$  then for all  $\varepsilon > 0$  there exists  $\delta > 0$  such that if  $|\mathbf{x} - \mathbf{y}| < \delta$  and  $\mathbf{y} \in D(\mathbf{f})$ , then  $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| < \varepsilon$ . Therefore, since  $\mathbf{x}$  is a limit point of  $D(\mathbf{f})$ ,  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{f}(\mathbf{x})$ . Suppose then that  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{f}(\mathbf{x})$  and let  $\varepsilon > 0$  be given. From the definition of the limit, there exists  $\delta > 0$  such that whenever  $\mathbf{y} \in D(\mathbf{f})$  and  $|\mathbf{y} - \mathbf{x}| < \delta$ , it follows  $|\mathbf{f}(\mathbf{y}) - \mathbf{f}(\mathbf{x})| < \varepsilon$ . Also,  $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x})| < \varepsilon$  and so if  $|\mathbf{x} - \mathbf{y}| < \delta$ , then  $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| < \varepsilon$ .

## 20.6 Exercises

1.  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y})$  was only defined in the case where  $\mathbf{x}$  was a limit point of  $D(\mathbf{f})$ . Why?
2. Suppose  $\mathbf{x}$  is defined to be a limit point of a set,  $A$  if and only if for all  $r > 0$ ,  $B(\mathbf{x}, r)$  contains a point of  $A$  different than  $\mathbf{x}$ . Show this is equivalent to the above definition of limit point.
3. Let  $\{\mathbf{x}_k\}_{k=1}^n$  be any finite set of points in  $\mathbb{R}^p$ . Show this set has no limit points.
4. Show that if  $\mathbf{f}, \mathbf{g} : U \rightarrow \mathbb{R}^3$  where  $U$  is a nonempty subset of  $\mathbb{R}^p$  and  $\mathbf{x}$  is a limit point of  $U$ , show if  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) = \mathbf{L}$  and  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{g}(\mathbf{y}) = \mathbf{M}$ , then  $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \mathbf{f}(\mathbf{y}) \times \mathbf{g}(\mathbf{y}) = \mathbf{L} \times \mathbf{M}$ .

5. Let

$$f(x, y) \equiv \begin{cases} \frac{xy}{x^2+y^2} & \text{if } (x, y) \neq (0, 0) \\ 0 & \text{if } (x, y) = (0, 0) \end{cases}.$$

Find  $\lim_{(x,y) \rightarrow (0,0)} f(x, y)$  if it exists. If it does not exist, tell why it does not exist.

6. Find  $\lim_{\mathbf{x} \rightarrow \mathbf{0}} \frac{\sin(|\mathbf{x}|)}{|\mathbf{x}|}$  and prove your answer from the definition of limit.
7. Suppose  $\mathbf{g}$  is a continuous vector valued function of one variable defined on  $[0, \infty)$ . Prove

$$\lim_{\mathbf{x} \rightarrow \mathbf{x}_0} \mathbf{g}(|\mathbf{x}|) = \mathbf{g}(|\mathbf{x}_0|).$$

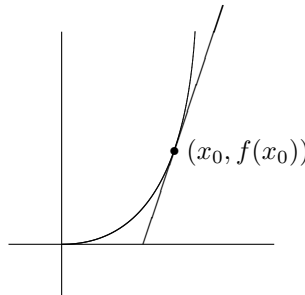
8. Give some examples of limit problems for functions of many variables which have limits and prove your assertions.

## 20.7 Approximation With A Tangent Plane

To begin with, suppose  $f$  is a function of one variable which has a derivative at the point  $x_0$ . Then the tangent line to the graph of the function,  $f$  at the point  $(x_0, f(x_0))$  is given by

$$y - f(x_0) = f'(x_0)(x - x_0) \quad (20.8)$$

and a picture of the situation is depicted below.



Thus from (20.8),

$$y = f(x_0) + f'(x_0)(x - x_0).$$

How close  $y$  is to  $f(x)$ ? Letting  $v = (x - x_0)$ ,

$$\begin{aligned} f(x_0 + v) - y &= \frac{f(x_0 + v) - f(x_0)}{v}(v) - f'(x_0)(v) \\ &= \left( \overbrace{\frac{f(x_0 + v) - f(x_0)}{v} - f'(x_0)}^{\rightarrow 0 \text{ as } v \rightarrow 0} \right) v. \end{aligned}$$

Now this shows

$$f(x_0 + v) - (f(x_0) + f'(x_0)v) = o(v)$$

where

$$\lim_{v \rightarrow 0} \frac{o(v)}{|v|} = 0.$$

Therefore, for small  $v$ ,  $f(x_0) + f'(x_0)v$  is a very good approximation to,  $f(x_0 + v)$ .

What of a function of two variables? Obviously such an approximation exists in one direction if the function has a directional derivative in the given direction, but what about letting  $\mathbf{v}$  be arbitrary and small? The idea similar to the above is an expression of the form

$$f((x_0, y_0) + (v_1, v_2)) = f(x_0, y_0) + av_1 + bv_2 + \overbrace{o(v_1, v_2)}^{=o(\mathbf{v})}$$

where

$$\lim_{|\mathbf{v}| \rightarrow 0} \frac{|o(\mathbf{v})|}{|\mathbf{v}|} = 0. \quad (20.9)$$

**Theorem 20.7.1** *Let  $U$  be an open subset of  $\mathbb{R}^2$  and suppose  $f : U \rightarrow \mathbb{R}$  has the property that the partial derivatives  $f_x$  and  $f_y$  exist for  $(x, y) \in U$  and are continuous at the point  $(x_0, y_0)$ . Then*

$$f((x_0, y_0) + (v_1, v_2)) = f(x_0, y_0) + \frac{\partial f}{\partial x}(x_0, y_0)v_1 + \frac{\partial f}{\partial y}(x_0, y_0)v_2 + o(\mathbf{v})$$

where  $o(\mathbf{v})$  satisfies (20.9).

**Proof:**

$$\begin{aligned} f((x_0, y_0) + (v_1, v_2)) &- \left( f(x_0, y_0) + \frac{\partial f}{\partial x}(x_0, y_0)v_1 + \frac{\partial f}{\partial y}(x_0, y_0)v_2 \right) \\ &= (f(x_0 + v_1, y_0 + v_2) - f(x_0, y_0)) - \left( \frac{\partial f}{\partial x}(x_0, y_0)v_1 + \frac{\partial f}{\partial y}(x_0, y_0)v_2 \right) \\ &= \left( \overbrace{f(x_0 + v_1, y_0 + v_2) - f(x_0, y_0 + v_2)}^{\text{changes only in first component}} + \overbrace{f(x_0, y_0 + v_2) - f(x_0, y_0)}^{\text{changes only in second component}} \right) \\ &\quad - \left( \frac{\partial f}{\partial x}(x_0, y_0)v_1 + \frac{\partial f}{\partial y}(x_0, y_0)v_2 \right) \end{aligned} \quad (20.10)$$

By the mean value theorem, there exist numbers  $s$  and  $t$  in  $[0, 1]$  such that this equals

$$= \left( \frac{\partial f}{\partial x}(x_0 + tv_1, y_0 + v_2)v_1 + \frac{\partial f}{\partial y}(x_0, y_0 + sv_2)v_2 \right)$$

$$\begin{aligned}
& - \left( \frac{\partial f}{\partial x}(x_0, y_0) v_1 + \frac{\partial f}{\partial y}(x_0, y_0) v_2 \right) \\
& = \left( \frac{\partial f}{\partial x}(x_0 + tv_1, y_0 + v_2) - \frac{\partial f}{\partial x}(x_0, y_0) \right) v_1 + \left( \frac{\partial f}{\partial y}(x_0, y_0 + sv_2) - \frac{\partial f}{\partial y}(x_0, y_0) \right) v_2
\end{aligned}$$

Therefore, letting  $o(\mathbf{v})$  denote the expression in (20.10), and noticing that  $|v_1|$  and  $|v_2|$  are both no larger than  $|\mathbf{v}|$ ,

$$|o(\mathbf{v})| \leq \left( \left| \frac{\partial f}{\partial x}(x_0 + tv_1, y_0 + v_2) - \frac{\partial f}{\partial x}(x_0, y_0) \right| + \left| \frac{\partial f}{\partial y}(x_0, y_0 + sv_2) - \frac{\partial f}{\partial y}(x_0, y_0) \right| \right) |\mathbf{v}|.$$

It follows

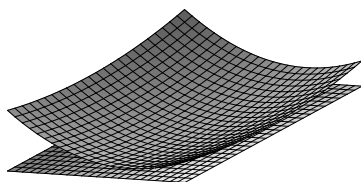
$$\frac{|o(\mathbf{v})|}{|\mathbf{v}|} \leq \left| \frac{\partial f}{\partial x}(x_0 + tv_1, y_0 + v_2) - \frac{\partial f}{\partial x}(x_0, y_0) \right| + \left| \frac{\partial f}{\partial y}(x_0, y_0 + sv_2) - \frac{\partial f}{\partial y}(x_0, y_0) \right|$$

Therefore,  $\lim_{\mathbf{v} \rightarrow \mathbf{0}} \frac{|o(\mathbf{v})|}{|\mathbf{v}|} = 0$  because of the assumption that  $f_x$  and  $f_y$  are continuous at the point  $(x_0, y_0)$  and this proves the theorem.

Writing this in a different form, letting  $\mathbf{v} \equiv (x - x_0, y - y_0)$ ,

$$f(x, y) = f(x_0, y_0) + \frac{\partial f}{\partial x}(x_0, y_0)(x - x_0) + \frac{\partial f}{\partial y}(x_0, y_0)(y - y_0) + o(\mathbf{v})$$

The right side of the above,  $f(x_0, y_0) + \frac{\partial f}{\partial x}(x_0, y_0)(x - x_0) + \frac{\partial f}{\partial y}(x_0, y_0)(y - y_0) = z$  is the equation of a plane tangent to the graph of  $f$ . The message of the above theorem is that the approximation is very good if both  $|x - x_0|$  and  $|y - y_0|$  are small.



Of course the above theorem has a generalization to any number of variables but for more than two it is impossible to draw a meaningful picture. Nevertheless, there is a meaningful theorem about approximating the function with one of a special form as above.

**Theorem 20.7.2** *Let  $U$  be an open subset of  $\mathbb{R}^p$  for  $p \geq 1$  and suppose  $f : U \rightarrow \mathbb{R}$  has the property that the partial derivatives  $f_{x_i}$  exist for all  $\mathbf{x} \in U$  and are continuous at the point  $\mathbf{x}_0 \in U$ . Then*

$$f(\mathbf{x}_0 + \mathbf{v}) = f(\mathbf{x}_0) + \sum_{i=1}^p \frac{\partial f}{\partial x_i}(\mathbf{x}_0) v_i + o(\mathbf{v}).$$

*That is,  $f$  is differentiable at  $\mathbf{x}_0$  and the derivative of  $f$  equals the linear transformation obtained by multiplying by the  $1 \times p$  matrix,*

$$\left( \frac{\partial f}{\partial x_1}(\mathbf{x}_0), \dots, \frac{\partial f}{\partial x_p}(\mathbf{x}_0) \right).$$

**Proof:** The proof is similar to the case of two variables. Letting  $\mathbf{v} = (v_1 \cdots, v_p)^T$ , denote by  $\theta_i \mathbf{v}$  the vector

$$(0, \dots, 0, v_i, v_{i+1}, \dots, v_p)^T$$



Thus  $\theta_0 \mathbf{v} = \mathbf{v}$ ,  $\theta_{p-1}(\mathbf{v}) = (0, \dots, 0, v_p)^T$ , and  $\theta_p \mathbf{v} = \mathbf{0}$ . Now

$$\begin{aligned} f(\mathbf{x}_0 + \mathbf{v}) &= \left( f(\mathbf{x}_0) + \sum_{i=1}^p \frac{\partial f}{\partial x_i}(\mathbf{x}_0) v_i \right) \\ &= \sum_{i=1}^p \left( \overbrace{f(\mathbf{x}_0 + \theta_{i-1} \mathbf{v}) - f(\mathbf{x}_0 + \theta_i \mathbf{v})}^{\text{changes only in the } i^{\text{th}} \text{ position}} \right) - \sum_{i=1}^p \frac{\partial f}{\partial x_i}(\mathbf{x}_0) v_i \end{aligned} \quad (20.11)$$

Now by the mean value theorem there exist numbers  $s_i \in (0, 1)$  such that the above expression equals

$$= \sum_{i=1}^p \frac{\partial f}{\partial x_i}(\mathbf{x}_0 + \theta_i \mathbf{v} + s_i v_i) v_i - \sum_{i=1}^p \frac{\partial f}{\partial x_i}(\mathbf{x}_0) v_i$$

and so letting  $o(\mathbf{v})$  equal the expression in (20.11),

$$\begin{aligned} |o(\mathbf{v})| &\leq \sum_{i=1}^p \left| \frac{\partial f}{\partial x_i}(\mathbf{x}_0 + \theta_i \mathbf{v} + s_i v_i) - \frac{\partial f}{\partial x_i}(\mathbf{x}_0) \right| |v_i| \\ &\leq \sum_{i=1}^p \left| \frac{\partial f}{\partial x_i}(\mathbf{x}_0 + \theta_i \mathbf{v} + s_i v_i) - \frac{\partial f}{\partial x_i}(\mathbf{x}_0) \right| |\mathbf{v}| \end{aligned}$$

and so

$$\lim_{\mathbf{v} \rightarrow \mathbf{0}} \frac{|o(\mathbf{v})|}{|\mathbf{v}|} \leq \lim_{\mathbf{v} \rightarrow \mathbf{0}} \sum_{i=1}^p \left| \frac{\partial f}{\partial x_i}(\mathbf{x}_0 + \theta_i \mathbf{v} + s_i v_i) - \frac{\partial f}{\partial x_i}(\mathbf{x}_0) \right| = 0$$

because of continuity of the  $f_{x_i}$  at  $\mathbf{x}_0$ . This proves the theorem.

As before, let  $\mathbf{x} - \mathbf{x}_0 = \mathbf{v}$  and write

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \sum_{i=1}^p \frac{\partial f}{\partial x_i}(\mathbf{x}_0) (x_i - x_{0i}) + o(\mathbf{v})$$

where  $o(\mathbf{v})$  satisfies (20.9). This is called the linear approximation of the function  $f$  at the point  $\mathbf{x}_0$ .

**Definition 20.7.3** Let  $\mathbf{f}: U \rightarrow \mathbb{R}^p$  where  $U$  is an open set in  $\mathbb{R}^n$ . Then  $\mathbf{f}$  is in  $C^1(U)$  if all its partial derivatives exist and are continuous on  $U$ . Sometimes people use  $C^1$  as an adjective and say  $\mathbf{f}$  is  $C^1$ .

**Example 20.7.4** Suppose  $f(x, y, z) = xy + z^2$ . Approximate  $f(x, y, z)$  for  $(x, y, z)$  near  $(1, 2, 3)$ .

Taking the partial derivatives of  $f$ ,  $f_x = y$ ,  $f_y = x$ ,  $f_z = 2z$ . Therefore,  $f_x(1, 2, 3) = 2$ ,  $f_y(1, 2, 3) = 1$ , and  $f_z(1, 2, 3) = 6$ . Therefore, the desired approximation is

$$\begin{aligned} &f(1, 2, 3) + 1(x - 1) + 2(y - 2) + 6(z - 3) \\ &= 11 + 1(x - 1) + 2(y - 2) + 6(z - 3) = -12 + x + 2y + 6z \end{aligned}$$

What about the case where  $\mathbf{f}$  has values in  $\mathbb{R}^q$  rather than  $\mathbb{R}$ ? Is there a similar theorem about linear approximations in this case also?

**Theorem 20.7.5** Let  $U$  be an open subset of  $\mathbb{R}^p$  for  $p \geq 1$  and suppose  $\mathbf{f} : U \rightarrow \mathbb{R}^q$  has the property that the partial derivatives  $\mathbf{f}_{x_i}$  exist for all  $\mathbf{x} \in U$  and are continuous at the point  $\mathbf{x}_0 \in U$ , then

$$\mathbf{f}(\mathbf{x}_0 + \mathbf{v}) = \mathbf{f}(\mathbf{x}_0) + \sum_{i=1}^p \frac{\partial \mathbf{f}}{\partial x_i}(\mathbf{x}_0) v_i + \mathbf{o}(\mathbf{v}) \quad (20.12)$$

where  $\mathbf{o}(\mathbf{v})$  satisfies (20.9).

**Proof:** Let  $\mathbf{f}(\mathbf{x}) \equiv (f_1(\mathbf{x}), \dots, f_q(\mathbf{x}))^T$ . From Theorem 20.7.2, the following holds for each  $k = 1, \dots, q$ .

$$f_k(\mathbf{x}) = f_k(\mathbf{x}_0) + \sum_{i=1}^p \frac{\partial f_k}{\partial x_i}(\mathbf{x}_0) v_i + o_k(\mathbf{v}).$$

Define  $\mathbf{o}(\mathbf{v}) \equiv (o_1(\mathbf{v}), \dots, o_q(\mathbf{v}))^T$ . Then (20.9) holds for  $\mathbf{o}(\mathbf{v})$  because it holds for each of the components of  $\mathbf{o}(\mathbf{v})$ . Also, the above equation is then equivalent to (20.12). Letting  $\mathbf{v} = \mathbf{x} - \mathbf{x}_0$ , this is equivalent to

$$\mathbf{f}(\mathbf{x}) = \mathbf{f}(\mathbf{x}_0) + \sum_{i=1}^p \frac{\partial \mathbf{f}}{\partial x_i}(\mathbf{x}_0) (x_i - x_{0i}) + \mathbf{o}(\mathbf{x} - \mathbf{x}_0). \quad (20.13)$$

When a formula like the above holds, the function must be continuous at  $\mathbf{x}_0$ . This is the content of the following important lemma.

**Lemma 20.7.6** Let  $\mathbf{f} : U \rightarrow \mathbb{R}^q$  where  $U$  is an open subset of  $\mathbb{R}^p$ . If (20.13) holds, then  $\mathbf{f}$  is continuous at  $\mathbf{x}_0$ . Furthermore, if  $C \geq \max \left\{ \left| \frac{\partial \mathbf{f}}{\partial x_i}(\mathbf{x}_0) \right|, i = 1, \dots, p \right\}$ , then whenever  $|\mathbf{x} - \mathbf{x}_0|$  is small enough,

$$|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}_0)| \leq (Cp + 1) |\mathbf{x} - \mathbf{x}_0| \quad (20.14)$$

**Proof:** Suppose (20.13) holds. Since  $\mathbf{o}(\mathbf{v})$  satisfies (20.9), there exists  $\delta_1 > 0$  such that if  $|\mathbf{x} - \mathbf{x}_0| < \delta_1$ , then  $|\mathbf{o}(\mathbf{x} - \mathbf{x}_0)| < |\mathbf{x} - \mathbf{x}_0|$ . But also, by the triangle inequality, Corollary 12.3.5 on Page 299,

$$\left| \sum_{i=1}^p \frac{\partial \mathbf{f}}{\partial x_i}(\mathbf{x}_0) (x_i - x_{0i}) \right| \leq C \sum_{i=1}^p |x_i - x_{0i}| \leq Cp |\mathbf{x} - \mathbf{x}_0|$$

Therefore, if  $|\mathbf{x} - \mathbf{x}_0| < \delta_1$ ,

$$\begin{aligned} |\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}_0)| &\leq \left| \sum_{i=1}^p \frac{\partial \mathbf{f}}{\partial x_i}(\mathbf{x}_0) (x_i - x_{0i}) \right| + |\mathbf{x} - \mathbf{x}_0| \\ &< (Cp + 1) |\mathbf{x} - \mathbf{x}_0| \end{aligned}$$

which verifies (20.14). Now letting  $\varepsilon > 0$  be given, let  $\delta = \min \left( \delta_1, \frac{\varepsilon}{Cp+1} \right)$ . Then for  $|\mathbf{x} - \mathbf{x}_0| < \delta$ ,

$$|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}_0)| < (Cp + 1) |\mathbf{x} - \mathbf{x}_0| < (Cp + 1) \frac{\varepsilon}{Cp + 1} = \varepsilon$$

showing  $\mathbf{f}$  is continuous at  $\mathbf{x}_0$ .

**Definition 20.7.7** A function,  $\mathbf{v} \rightarrow \mathbf{g}(\mathbf{v})$  is  $\mathbf{o}(\mathbf{v})$  if (20.9) holds for  $\mathbf{g}(\mathbf{v})$ . That is,

$$\lim_{\mathbf{v} \rightarrow \mathbf{0}} \frac{|\mathbf{g}(\mathbf{v})|}{|\mathbf{v}|} = 0.$$

Also, the symbol,  $\mathbf{o}(\mathbf{v})$ , means  $\mathbf{o}(\cdot)$  is a function which satisfies (20.9). Thus it is permissible to write things like  $\mathbf{o}(\mathbf{v}) = 10\mathbf{o}(\mathbf{v})$  or  $\mathbf{o}(\mathbf{v}) = \mathbf{o}(\mathbf{v}) + \mathbf{o}(\mathbf{v})$  because if  $\mathbf{o}(\mathbf{v})$  satisfies (20.9), then so does  $\mathbf{o}(\mathbf{v}) + \mathbf{o}(\mathbf{v})$  and  $10\mathbf{o}(\mathbf{v})$ .

This is both imprecise and useful because it neglects everything which is not important and directs our attention to that which really matters, namely (20.9). I find it helpful to think of  $\mathbf{o}(\mathbf{v})$  as an adjective rather than a precise definition.

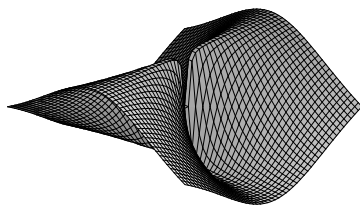
Something more than existence of the partial derivatives at a point is necessary in order to accomplish such a linear approximation. Here is an example.

$$f(x, y) \equiv \begin{cases} \frac{xy}{x^2+y^2} & \text{if } (x, y) \neq (0, 0) \\ 0 & \text{if } (x, y) = (0, 0) \end{cases}$$

Then

$$f_x(0, 0) \equiv \lim_{t \rightarrow 0} \frac{f(t, 0) - f(0, 0)}{t} = \lim_{t \rightarrow 0} \frac{0}{t} = 0.$$

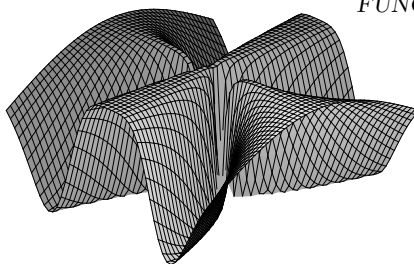
Similarly,  $f_y(0, 0) = 0$  so the partial derivatives do indeed exist. However, this function is not even continuous at the point  $(0, 0)$  because for every  $\delta > 0$ ,  $f(\delta/2, \delta/2) = \frac{1}{2}$  while  $f(0, \delta/2) = 0$ . Thus  $B((0, 0), \delta)$  contains points where the value of the function equals  $1/2$  and also points where the value of the function equals 0. Consequently no limit can exist. (Why?) It follows from Theorem 20.5.6 on Page 477 that  $f$  is not continuous at  $(0, 0)$ . Therefore, from the above discussion, it is impossible to give the sort of approximation for  $f$  described in Theorem 20.7.1 at the point  $(0, 0)$ . You might find the picture of this function amusing.



You see the way it is pinched near the center. I think you should notice that from the above discussion you see there is a problem and so it can be seen in the picture because it was expected. However, it is likely the picture would never have revealed the truth. This illustrates the limitations of abandoning careful logical reasoning and definitions in favor of pictures. I realize some of the above theorems and definitions are hard but they are the only route to correct understanding of multivariable calculus.

Here is another picture. This time it is of the function,

$$f(x, y) \equiv \begin{cases} \left( \frac{x^2 - y^4}{x^2 + y^4} \right)^2 & \text{if } (x, y) \neq (0, 0) \\ 1 & \text{if } (x, y) = (0, 0) \end{cases} \quad (20.15)$$



This is a very interesting example which you will examine in the exercises.

## 20.8 Exercises

1. Show the function in (20.15) is not continuous at  $(0,0)$  and yet it has directional derivatives in every direction at  $(0,0)$ .
2. Find a linear approximation to the function  $f(x, y, z) = x\sqrt[3]{y}z^2$  near the point  $(1, 8, 1)$ . Use this approximation to estimate  $f(1.1, 7.9, 1.1)$  and then use a calculator to make a comparison.
3. In 1 Kings chapter 7 there is a description of the molten sea in Solomon's temple. "And he made a molten sea, ten cubits from the one brim to the other: it was round all about, and his height was five cubits: and a line of thirty cubits did compass it round about..... And it was an hand breadth thick.... It contained two thousand baths." Lets assume a hand breadth is 8 inches and a cubit is 1.5 feet. What was the approximate volume in cubic feet of the brass used in the molten sea. How many baths of brass were used?

## 20.9 Differentiation And The Chain Rule

### 20.9.1 The Chain Rule

As in the case of a function of one variable, it is important to consider the derivative of a composition of two functions. A function,  $\mathbf{g} : U \rightarrow \mathbb{R}^p$  has a good linear approximation at  $\mathbf{x} \in U$  if

$$\mathbf{g}(\mathbf{x} + \mathbf{v}) = \mathbf{g}(\mathbf{x}) + \sum_{i=1}^n \frac{\partial \mathbf{g}(\mathbf{x})}{\partial x_i} v_i + \mathbf{o}(\mathbf{v}). \quad (20.16)$$

**Lemma 20.9.1** *Let  $\mathbf{g} : U \rightarrow \mathbb{R}^p$  where  $U$  is an open set in  $\mathbb{R}^n$  and suppose  $\mathbf{g}$  has a good linear approximation at  $\mathbf{x} \in U$ . Then  $\mathbf{o}(\mathbf{g}(\mathbf{x} + \mathbf{v}) - \mathbf{g}(\mathbf{x})) = \mathbf{o}(\mathbf{v})$ .*

**Proof:** It is necessary to show

$$\lim_{\mathbf{v} \rightarrow \mathbf{0}} \frac{|\mathbf{o}(\mathbf{g}(\mathbf{x} + \mathbf{v}) - \mathbf{g}(\mathbf{x}))|}{|\mathbf{v}|} = 0. \quad (20.17)$$

From (20.16), and Lemma 20.7.6, there exists  $\delta > 0$  such that if  $|\mathbf{v}| < \delta$ , then

$$|\mathbf{g}(\mathbf{x} + \mathbf{v}) - \mathbf{g}(\mathbf{x})| \leq (Cn + 1) |\mathbf{v}|. \quad (20.18)$$

Now let  $\varepsilon > 0$  be given. There exists  $\eta > 0$  such that if  $|\mathbf{g}(\mathbf{x} + \mathbf{v}) - \mathbf{g}(\mathbf{x})| < \eta$ , then

$$|\mathbf{o}(\mathbf{g}(\mathbf{x} + \mathbf{v}) - \mathbf{g}(\mathbf{x}))| < \left( \frac{\varepsilon}{Cn + 1} \right) |\mathbf{g}(\mathbf{x} + \mathbf{v}) - \mathbf{g}(\mathbf{x})| \quad (20.19)$$

Let  $|\mathbf{v}| < \min\left(\delta, \frac{\eta}{Cn+1}\right)$ . For such  $\mathbf{v}$ ,  $|\mathbf{g}(\mathbf{x} + \mathbf{v}) - \mathbf{g}(\mathbf{x})| \leq \eta$ , which implies

$$\begin{aligned} |\mathbf{o}(\mathbf{g}(\mathbf{x} + \mathbf{v}) - \mathbf{g}(\mathbf{x}))| &< \left(\frac{\varepsilon}{Cn+1}\right) |\mathbf{g}(\mathbf{x} + \mathbf{v}) - \mathbf{g}(\mathbf{x})| \\ &< \left(\frac{\varepsilon}{Cn+1}\right) (Cn+1) |\mathbf{v}| \end{aligned}$$

and so

$$\frac{|\mathbf{o}(\mathbf{g}(\mathbf{x} + \mathbf{v}) - \mathbf{g}(\mathbf{x}))|}{|\mathbf{v}|} < \varepsilon$$

which establishes (20.17). This proves the lemma.

Recall the notation  $\mathbf{f} \circ \mathbf{g}(\mathbf{x}) \equiv \mathbf{f}(\mathbf{g}(\mathbf{x}))$ . Thus  $\mathbf{f} \circ \mathbf{g}$  is the name of a function and this function is defined by what was just written. The last assertion of the following theorem is known as the chain rule.

**Theorem 20.9.2** *Let  $U$  be an open set in  $\mathbb{R}^n$ , let  $V$  be an open set in  $\mathbb{R}^p$ , let  $\mathbf{g} : U \rightarrow \mathbb{R}^p$  be such that  $\mathbf{g}(U) \subseteq V$ , and let  $\mathbf{f} : V \rightarrow \mathbb{R}^q$ . Suppose  $\mathbf{g}$  has a good linear approximation at  $\mathbf{x} \in U$  and that  $\mathbf{f}$  has a good linear approximation at  $\mathbf{g}(\mathbf{x})$ . Then  $\mathbf{f} \circ \mathbf{g}$  has a good linear approximation at  $\mathbf{x}$  and furthermore,*

$$\mathbf{f}(\mathbf{g}(\mathbf{x} + \mathbf{v})) = \mathbf{f}(\mathbf{g}(\mathbf{x})) + \sum_{j=1}^n \left( \sum_{i=1}^p \frac{\partial \mathbf{f}(\mathbf{g}(\mathbf{x}))}{\partial y_i} \frac{\partial g_i(\mathbf{x})}{\partial x_j} \right) v_j + \mathbf{o}(\mathbf{v}). \quad (20.20)$$

and

$$\frac{\partial (\mathbf{f} \circ \mathbf{g})(\mathbf{x})}{\partial x_j} = \sum_{i=1}^p \frac{\partial \mathbf{f}(\mathbf{g}(\mathbf{x}))}{\partial y_i} \frac{\partial g_i(\mathbf{x})}{\partial x_j}. \quad (20.21)$$

**Proof:** From the assumption that  $\mathbf{f}$  has a good linear approximation at  $\mathbf{g}(\mathbf{x})$ ,

$$\mathbf{f}(\mathbf{g}(\mathbf{x} + \mathbf{v})) = \mathbf{f}(\mathbf{g}(\mathbf{x})) + \sum_{i=1}^p \frac{\partial \mathbf{f}(\mathbf{g}(\mathbf{x}))}{\partial y_i} (g_i(\mathbf{x} + \mathbf{v}) - g_i(\mathbf{x})) + \mathbf{o}(\mathbf{g}(\mathbf{x} + \mathbf{v}) - \mathbf{g}(\mathbf{x}))$$

which by Lemma 20.9.1 equals

$$\mathbf{f}(\mathbf{g}(\mathbf{x} + \mathbf{v})) = \mathbf{f}(\mathbf{g}(\mathbf{x})) + \sum_{i=1}^p \frac{\partial \mathbf{f}(\mathbf{g}(\mathbf{x}))}{\partial y_i} (g_i(\mathbf{x} + \mathbf{v}) - g_i(\mathbf{x})) + \mathbf{o}(\mathbf{v}).$$

Now since  $\mathbf{g}$  has a good linear approximation at  $\mathbf{x}$ , the above becomes

$$\begin{aligned} \mathbf{f}(\mathbf{g}(\mathbf{x} + \mathbf{v})) &= \mathbf{f}(\mathbf{g}(\mathbf{x})) + \sum_{i=1}^p \frac{\partial \mathbf{f}(\mathbf{g}(\mathbf{x}))}{\partial y_i} \left( \sum_{j=1}^n \frac{\partial g_i(\mathbf{x})}{\partial x_j} v_j + \mathbf{o}(\mathbf{v}) \right) + \mathbf{o}(\mathbf{v}) \\ &= \mathbf{f}(\mathbf{g}(\mathbf{x})) + \sum_{i=1}^p \frac{\partial \mathbf{f}(\mathbf{g}(\mathbf{x}))}{\partial y_i} \left( \sum_{j=1}^n \frac{\partial g_i(\mathbf{x})}{\partial x_j} v_j \right) + \sum_{i=1}^p \frac{\partial \mathbf{f}(\mathbf{g}(\mathbf{x}))}{\partial y_i} \mathbf{o}(\mathbf{v}) + \mathbf{o}(\mathbf{v}) \\ &= \mathbf{f}(\mathbf{g}(\mathbf{x})) + \sum_{j=1}^n \left( \sum_{i=1}^p \frac{\partial \mathbf{f}(\mathbf{g}(\mathbf{x}))}{\partial y_i} \frac{\partial g_i(\mathbf{x})}{\partial x_j} \right) v_j + \mathbf{o}(\mathbf{v}) \end{aligned}$$

because  $\sum_{i=1}^p \frac{\partial \mathbf{f}(\mathbf{g}(\mathbf{x}))}{\partial y_i} \mathbf{o}(\mathbf{v}) + \mathbf{o}(\mathbf{v}) = \mathbf{o}(\mathbf{v})$ . This establishes (20.20). It only remains to verify (20.21). Using the formula,

$$\begin{aligned} \frac{\partial (\mathbf{f} \circ \mathbf{g})(\mathbf{x})}{\partial x_k} &= \lim_{t \rightarrow 0} \frac{\mathbf{f}(\mathbf{g}(\mathbf{x} + t\mathbf{e}_k)) - \mathbf{f}(\mathbf{g}(\mathbf{x}))}{t} \\ &= \lim_{t \rightarrow 0} \frac{1}{t} \left[ \left( \sum_{i=1}^p \frac{\partial \mathbf{f}(\mathbf{g}(\mathbf{x}))}{\partial y_i} \frac{\partial g_i(\mathbf{x})}{\partial x_k} \right) t + \mathbf{o}(t\mathbf{e}_k) \right]. \end{aligned}$$

This follows from observing that for  $\mathbf{v} = t\mathbf{e}_k$ ,  $v_i = 0$  unless  $i = k$  in which case  $v_k = t$ . Now

$$\lim_{t \rightarrow 0} \frac{|\mathbf{o}(t\mathbf{e}_k)|}{t} = 0$$

and so, taking the limit yields

$$\frac{\partial (\mathbf{f} \circ \mathbf{g})(\mathbf{x})}{\partial x_k} = \sum_{i=1}^p \frac{\partial \mathbf{f}(\mathbf{g}(\mathbf{x}))}{\partial y_i} \frac{\partial g_i(\mathbf{x})}{\partial x_k}$$

as claimed.

There is an easy way to remember this in terms of the repeated index summation convention presented earlier. Let  $\mathbf{y} = \mathbf{g}(\mathbf{x})$  and  $\mathbf{z} = \mathbf{f}(\mathbf{y})$ . Then the above says

$$\frac{\partial \mathbf{z}}{\partial y_i} \frac{\partial y_i}{\partial x_k} = \frac{\partial \mathbf{z}}{\partial x_k}.$$

Remember there is a sum on the repeated index.

**Example 20.9.3** Let  $f(u, v) = \sin(uv)$  and let  $u(x, y, t) = t \sin x + \cos y$  and  $v(x, y, t, s) = s \tan x + y^2 + ts$ . Letting  $z = f(u, v)$  where  $u, v$  are as just described, find  $\frac{\partial z}{\partial t}$  and  $\frac{\partial z}{\partial x}$ .

From the above,

$$\frac{\partial z}{\partial t} = \frac{\partial z}{\partial u} \frac{\partial u}{\partial t} + \frac{\partial z}{\partial v} \frac{\partial v}{\partial t} = v \cos(uv) \sin(x) + us \cos(uv).$$

Also,

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial u} \frac{\partial u}{\partial x} + \frac{\partial z}{\partial v} \frac{\partial v}{\partial x} = v \cos(uv) t \cos(x) + us \sec^2(x) \cos(uv).$$

Clearly you can continue in this way taking partial derivatives with respect to any of the other variables.

**Example 20.9.4** Recall spherical coordinates are given by

$$x = \rho \sin \phi \cos \theta, \quad y = \rho \sin \phi \sin \theta, \quad z = \rho \cos \phi.$$

If an object moves in three dimensions, describe its acceleration in terms of spherical coordinates and the vectors,

$$\mathbf{e}_\rho = (\sin \phi \cos \theta, \sin \phi \sin \theta, \cos \phi)^T,$$

$$\mathbf{e}_\theta = (-\rho \sin \phi \sin \theta, \rho \sin \phi \cos \theta, 0)^T,$$

and

$$\mathbf{e}_\phi = (\rho \cos \phi \cos \theta, \rho \cos \phi \sin \theta, -\rho \sin \phi)^T.$$

Why these vectors? Note how they were obtained. Let

$$\mathbf{r}(\rho, \theta, \phi) = (\rho \sin \phi \cos \theta, \rho \sin \phi \sin \theta, \rho \cos \phi)^T$$

and fix  $\phi$  and  $\theta$ , letting only  $\rho$  change, this gives a curve in the direction of increasing  $\rho$ . Thus it is a vector which points away from the origin. Letting only  $\phi$  change and fixing  $\theta$  and  $\rho$ , this gives a vector which is tangent to the sphere of radius  $\rho$  and points South. Similarly, letting  $\theta$  change and fixing the other two gives a vector which points East and is tangent to the sphere of radius  $\rho$ . It is thought by most people that we live on a large sphere. The giant turtle carrying the flat earth on its back as it proceeds through the universe is not usually accepted as a correct model. Given we live on a sphere, what directions would be most meaningful? Wouldn't it be the directions of the vectors just described?

Let  $\mathbf{r}(t)$  denote the position vector of the object from the origin. Thus

$$\mathbf{r}(t) = \rho(t) \mathbf{e}_\rho(t) = \left( (x(t), y(t), z(t))^T \right)$$

Now this implies the velocity is

$$\mathbf{r}'(t) = \rho'(t) \mathbf{e}_\rho(t) + \rho(t) (\mathbf{e}_\rho(t))'. \quad (20.22)$$

You see,  $\mathbf{e}_\rho = \mathbf{e}_\rho(\rho, \theta, \phi)$  where each of these variables is a function of  $t$ .

$$\frac{\partial \mathbf{e}_\rho}{\partial \phi} = (\cos \phi \cos \theta, \cos \phi \sin \theta, -\sin \phi)^T = \frac{1}{\rho} \mathbf{e}_\phi,$$

$$\frac{\partial \mathbf{e}_\rho}{\partial \theta} = (-\sin \phi \sin \theta, \sin \phi \cos \theta, 0)^T = \frac{1}{\rho} \mathbf{e}_\theta,$$

and

$$\frac{\partial \mathbf{e}_\rho}{\partial \rho} = 0.$$

Therefore, by the chain rule,

$$\begin{aligned} \frac{d\mathbf{e}_\rho}{dt} &= \frac{\partial \mathbf{e}_\rho}{\partial \phi} \frac{d\phi}{dt} + \frac{\partial \mathbf{e}_\rho}{\partial \theta} \frac{d\theta}{dt} \\ &= \frac{1}{\rho} \frac{d\phi}{dt} \mathbf{e}_\phi + \frac{1}{\rho} \frac{d\theta}{dt} \mathbf{e}_\theta. \end{aligned}$$

By (20.22),

$$\mathbf{r}' = \rho' \mathbf{e}_\rho + \frac{d\phi}{dt} \mathbf{e}_\phi + \frac{d\theta}{dt} \mathbf{e}_\theta. \quad (20.23)$$

Now things get interesting. This must be differentiated with respect to  $t$ . To do so,

$$\frac{\partial \mathbf{e}_\theta}{\partial \theta} = (-\rho \sin \phi \cos \theta, -\rho \sin \phi \sin \theta, 0)^T = ?$$

where it is desired to find  $a, b, c$  such that  $? = a\mathbf{e}_\theta + b\mathbf{e}_\phi + c\mathbf{e}_\rho$ . Thus

$$\begin{pmatrix} -\rho \sin \phi \sin \theta & \rho \cos \phi \cos \theta & \sin \phi \cos \theta \\ \rho \sin \phi \cos \theta & \rho \cos \phi \sin \theta & \sin \phi \sin \theta \\ 0 & -\rho \sin \phi & \cos \phi \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} -\rho \sin \phi \cos \theta \\ -\rho \sin \phi \sin \theta \\ 0 \end{pmatrix}$$

Using Cramer's rule, the solution is  $a = 0$ ,  $b = -\cos \phi \sin \phi$ , and  $c = -\rho \sin^2 \phi$ . Thus

$$\begin{aligned} \frac{\partial \mathbf{e}_\theta}{\partial \theta} &= (-\rho \sin \phi \cos \theta, -\rho \sin \phi \sin \theta, 0)^T \\ &= (-\cos \phi \sin \phi) \mathbf{e}_\phi + (-\rho \sin^2 \phi) \mathbf{e}_\rho. \end{aligned}$$

Also,

$$\frac{\partial \mathbf{e}_\theta}{\partial \phi} = (-\rho \cos \phi \sin \theta, \rho \cos \phi \cos \theta, 0)^T = (\cot \phi) \mathbf{e}_\theta$$

and

$$\frac{\partial \mathbf{e}_\theta}{\partial \rho} = (-\sin \phi \sin \theta, \sin \phi \cos \theta, 0)^T = \frac{1}{\rho} \mathbf{e}_\theta.$$

Now in (20.23) it is also necessary to consider  $\mathbf{e}_\phi$ .

$$\frac{\partial \mathbf{e}_\phi}{\partial \phi} = (-\rho \sin \phi \cos \theta, -\rho \sin \phi \sin \theta, -\rho \cos \phi)^T = -\rho \mathbf{e}_\rho$$

$$\begin{aligned} \frac{\partial \mathbf{e}_\phi}{\partial \theta} &= (-\rho \cos \phi \sin \theta, \rho \cos \phi \cos \theta, 0)^T \\ &= (\cot \phi) \mathbf{e}_\theta \end{aligned}$$

and finally,

$$\frac{\partial \mathbf{e}_\phi}{\partial \rho} = (\cos \phi \cos \theta, \cos \phi \sin \theta, -\sin \phi)^T = \frac{1}{\rho} \mathbf{e}_\phi.$$

With these formulas for various partial derivatives, the chain rule is used to obtain  $\mathbf{r}''$  which will yield a formula for the acceleration in terms of the spherical coordinates and these special vectors. By the chain rule,

$$\begin{aligned} \frac{d}{dt}(\mathbf{e}_\rho) &= \frac{\partial \mathbf{e}_\rho}{\partial \theta} \theta' + \frac{\partial \mathbf{e}_\rho}{\partial \phi} \phi' + \frac{\partial \mathbf{e}_\rho}{\partial \rho} \rho' \\ &= \frac{\theta'}{\rho} \mathbf{e}_\theta + \frac{\phi'}{\rho} \mathbf{e}_\phi \end{aligned}$$

$$\begin{aligned} \frac{d}{dt}(\mathbf{e}_\theta) &= \frac{\partial \mathbf{e}_\theta}{\partial \theta} \theta' + \frac{\partial \mathbf{e}_\theta}{\partial \phi} \phi' + \frac{\partial \mathbf{e}_\theta}{\partial \rho} \rho' \\ &= \theta' ((-\cos \phi \sin \phi) \mathbf{e}_\phi + (-\rho \sin^2 \phi) \mathbf{e}_\rho) + \phi' (\cot \phi) \mathbf{e}_\theta + \frac{\rho'}{\rho} \mathbf{e}_\theta \end{aligned}$$

$$\begin{aligned} \frac{d}{dt}(\mathbf{e}_\phi) &= \frac{\partial \mathbf{e}_\phi}{\partial \theta} \theta' + \frac{\partial \mathbf{e}_\phi}{\partial \phi} \phi' + \frac{\partial \mathbf{e}_\phi}{\partial \rho} \rho' \\ &= (\theta' \cot \phi) \mathbf{e}_\theta + \phi' (-\rho \mathbf{e}_\rho) + \left( \frac{\rho'}{\rho} \mathbf{e}_\phi \right) \end{aligned}$$

By (20.23),

$$\mathbf{r}'' = \rho'' \mathbf{e}_\rho + \phi'' \mathbf{e}_\phi + \theta'' \mathbf{e}_\theta + \rho' (\mathbf{e}_\rho)' + \phi' (\mathbf{e}_\phi)' + \theta' (\mathbf{e}_\theta)'$$

and from the above, this equals

$$\begin{aligned} &\rho'' \mathbf{e}_\rho + \phi'' \mathbf{e}_\phi + \theta'' \mathbf{e}_\theta + \rho' \left( \frac{\theta'}{\rho} \mathbf{e}_\theta + \frac{\phi'}{\rho} \mathbf{e}_\phi \right) + \\ &\phi' \left( (\theta' \cot \phi) \mathbf{e}_\theta + \phi' (-\rho \mathbf{e}_\rho) + \left( \frac{\rho'}{\rho} \mathbf{e}_\phi \right) \right) + \\ &\theta' \left( \theta' ((-\cos \phi \sin \phi) \mathbf{e}_\phi + (-\rho \sin^2 \phi) \mathbf{e}_\rho) + \phi' (\cot \phi) \mathbf{e}_\theta + \frac{\rho'}{\rho} \mathbf{e}_\theta \right) \end{aligned}$$



and now all that remains is to collect the terms. Thus  $\mathbf{r}''$  equals

$$\begin{aligned}\mathbf{r}'' &= \left( \rho'' - \rho (\phi')^2 - \rho (\theta')^2 \sin^2(\phi) \right) \mathbf{e}_\rho + \left( \phi'' + \frac{2\rho'\phi'}{\rho} - (\theta')^2 \cos\phi \sin\phi \right) \mathbf{e}_\phi + \\ &\quad + \left( \theta'' + \frac{2\theta'\rho'}{\rho} + 2\phi'\theta' \cot(\phi) \right) \mathbf{e}_\theta.\end{aligned}$$

and this gives the acceleration in spherical coordinates. Note the prominent role played by the chain rule. All of the above is done in books on mechanics for general curvilinear coordinate systems and in the more general context, special theorems are developed which make things go much faster but these theorems are all exercises in the chain rule.

As an example of how this could be used, consider a rocket. Suppose for simplicity that it experiences a force only in the direction of  $\mathbf{e}_\rho$ , directly away from the earth. Of course this force produces a corresponding acceleration which can be computed as a function of time. As the fuel is burned, the rocket becomes less massive and so the acceleration will be an increasing function of  $t$ . However, this would be a known function, say  $a(t)$ . Suppose you wanted to know the latitude and longitude of the rocket as a function of time. (There is no reason to think these will stay the same.) Then all that would be required would be to solve the system of differential equations<sup>3</sup>,

$$\begin{aligned}\rho'' - \rho (\phi')^2 - \rho (\theta')^2 \sin^2(\phi) &= a(t), \\ \phi'' + \frac{2\rho'\phi'}{\rho} - (\theta')^2 \cos\phi \sin\phi &= 0, \\ \theta'' + \frac{2\theta'\rho'}{\rho} + 2\phi'\theta' \cot(\phi) &= 0\end{aligned}$$

along with initial conditions,  $\rho(0) = \rho_0$  (the distance from the launch site to the center of the earth.),  $\rho'(0) = \rho_1$  (the initial vertical component of velocity of the rocket, probably 0.) and then initial conditions for  $\phi, \phi', \theta, \theta'$ . The initial value problems could then be solved numerically and you would know the distance from the center of the earth as a function of  $t$  along with  $\theta$  and  $\phi$ . Thus you could predict where the booster shells would fall to earth so you would know where to look for them. Of course there are many variations of this. You might want to specify forces in the  $\mathbf{e}_\theta$  and  $\mathbf{e}_\phi$  direction as well and attempt to control the position of the rocket or rather its payload. The point is that if you are interested in doing all this in terms of  $\phi, \theta$ , and  $\rho$ , the above shows how to do it systematically and you see it is all an exercise in using the chain rule. More could be said here involving moving coordinate systems and the Coriolis force.

### 20.9.2 Differentiation And The Derivative

By now, you may have suspected there should be a definition of something which will be equivalent to the assertion that a function has a good linear approximation. This is in fact the case.

**Definition 20.9.5** Let  $\mathbf{f} : U \rightarrow \mathbb{R}^p$  where  $U$  is an open set in  $\mathbb{R}^n$  for  $n, p \geq 1$  and let  $\mathbf{x} \in U$  be given. Then  $\mathbf{f}$  is defined to be differentiable at  $\mathbf{x} \in U$  if and only if  $\mathbf{f}$  has a good linear approximation at  $\mathbf{x}$ .

<sup>3</sup>You won't be able to find the solution to equations like these in terms of simple functions. However, they can be solved numerically. This means you determine the value of the various variables for various values of  $t$  without finding a neat formula involving known functions for the solution. This sort of computation is usually done by a computer.

Of course an obvious question arises. What if  $n, p = 1$ . Is this new definition equivalent to the old one in this case? In the case where  $n, p = 1$ , suppose  $f'(x)$  exists. Then

$$\lim_{|v| \rightarrow 0} \frac{|f(x+v) - f(x) - f'(x)v|}{|v|} = 0$$

because

$$\frac{|f(x+v) - f(x) - f'(x)v|}{|v|} = \left| \frac{f(x+v) - f(x)}{v} - f'(x) \right|$$

and  $f'(x)$  is defined as the  $\lim_{v \rightarrow 0} \frac{f(x+v) - f(x)}{v}$ . It follows that

$$f(x+v) - f(x) - f'(x)v = o(v).$$

Therefore, if  $f$  is differentiable in the old way, it is differentiable in the new way just described.

Now suppose  $f$  is differentiable in the new way in terms of having a good linear approximation. Then

$$\begin{aligned} \left| \frac{f(x+v) - f(x)}{v} - f'(x) \right| &= \left| \frac{f'(x)v + o(v)}{v} - f'(x) \right| \\ &= \left| \frac{o(v)}{v} \right| \end{aligned}$$

and this last expression converges to 0 by the definition of what it means to be  $o(v)$ . Therefore,  $f'(x)$  equals the limit of the difference quotient as before and this shows that if  $f$  is differentiable in the new way, then it is differentiable in the old way.

Why bother with a new definition if there is an old one which is already perfectly good? It is because in the case of a function of  $n$  variables, the old definition makes absolutely no sense because you can't divide by a vector. This is the whole reason for the new way of looking at things.

Now that a definition of what it means for a function to be differentiable has been given, an obvious question is: What is the derivative? In the case of a scalar function of one variable,

$$f(x+v) = f(x) + f'(x)v + o(v), \quad (20.24)$$

where  $v \in \mathbb{R}$ . Thus  $v \rightarrow f'(x)v$  is a linear mapping from  $\mathbb{R}$  to  $\mathbb{R}$ , and  $f'(x)$  can be thought of as this linear mapping. Till now, it has been thought of as a number and this is what it is, but multiplication by it yields a linear transformation from  $\mathbb{R}$  to  $\mathbb{R}$ . In the more general setting described above, the derivative is the same thing, a linear transformation.

Let  $\mathbf{f} : U \rightarrow \mathbb{R}^q$  where  $U$  is an open subset of  $\mathbb{R}^p$  and  $\mathbf{f}$  is differentiable. Recall this implies

$$\mathbf{f}(\mathbf{x} + \mathbf{v}) = \mathbf{f}(\mathbf{x}) + \sum_{j=1}^p \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_j} v_j + \mathbf{o}(\mathbf{v}).$$

Taking the  $i^{th}$  coordinate of the above equation yields

$$f_i(\mathbf{x} + \mathbf{v}) = f_i(\mathbf{x}) + \sum_{j=1}^p \frac{\partial f_i(\mathbf{x})}{\partial x_j} v_j + o(\mathbf{v})$$

and it follows that the term with a sum is nothing more than the  $i^{th}$  component of  $J(\mathbf{x}) \mathbf{v}$  where  $J(\mathbf{x})$  is the  $q \times p$  matrix,

$$\begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_p} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_p} \\ \vdots & & \ddots & \\ \frac{\partial f_q}{\partial x_1} & \frac{\partial f_q}{\partial x_2} & \cdots & \frac{\partial f_q}{\partial x_p} \end{pmatrix}$$

and that

$$\mathbf{f}(\mathbf{x} + \mathbf{v}) = \mathbf{f}(\mathbf{x}) + J(\mathbf{x}) \mathbf{v} + \mathbf{o}(\mathbf{v}). \quad (20.25)$$

The linear transformation which results by multiplication by this  $q \times p$  matrix is known as the derivative. This linear transformation is also denoted by  $D\mathbf{f}(\mathbf{x})$  or sometimes  $\mathbf{f}'(\mathbf{x})$  to correspond completely with the earlier notation. Thus

$$J(\mathbf{x}) \mathbf{v} = D\mathbf{f}(\mathbf{x}) \mathbf{v} = \mathbf{f}'(\mathbf{x}) \mathbf{v},$$

where the first of these terms means matrix multiplication in which  $\mathbf{v} \in \mathbb{R}^p$  and the second two denote a linear transformation acting on the vector,  $\mathbf{v}$ . There is no need to distinguish these subtle issues at this time. Just think of them as different notations for the derivative. There is also a notion of partial derivative in this context.

**Definition 20.9.6** Suppose  $U$  is an open set in  $\mathbb{R}^n$  and  $V$  is an open set in  $\mathbb{R}^p$ . Also suppose  $\mathbf{f} : U \times V \rightarrow \mathbb{R}^q$  is in  $C^1(U \times V)$ . Define  $D_1\mathbf{f}(\mathbf{x}, \mathbf{y})$  and  $D_2\mathbf{f}(\mathbf{x}, \mathbf{y})$  as follows.  $D_1\mathbf{f}(\mathbf{x}, \mathbf{y})$  is a linear transformation mapping  $\mathbb{R}^n$  to  $\mathbb{R}^q$  which satisfies

$$\mathbf{f}(\mathbf{x} + \mathbf{v}, \mathbf{y}) = \mathbf{f}(\mathbf{x}, \mathbf{y}) + D_1\mathbf{f}(\mathbf{x}, \mathbf{y}) \mathbf{v} + \mathbf{o}(\mathbf{v}).$$

(See Problem 1) In other words, fix  $\mathbf{y}$  and take the derivative of the resulting function of  $\mathbf{x}$ .  $D_2\mathbf{f}(\mathbf{x}, \mathbf{y})$  is defined similarly as the derivative of what is obtained when  $\mathbf{x}$  is fixed and the variable of interest is  $\mathbf{y}$ . Generalization to more factors follows the same pattern.

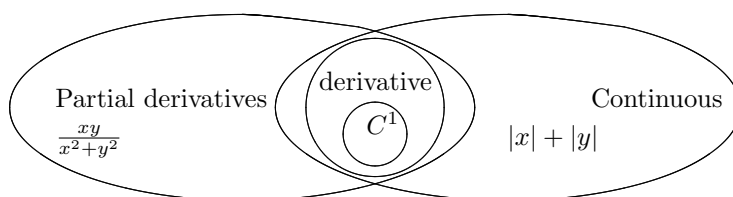
In terms of matrices of partial derivatives,

$$D_1\mathbf{f}(\mathbf{x}, \mathbf{y}) = \begin{pmatrix} \frac{\partial f_1(\mathbf{x}, \mathbf{y})}{\partial x_1} & \frac{\partial f_1(\mathbf{x}, \mathbf{y})}{\partial x_2} & \cdots & \frac{\partial f_1(\mathbf{x}, \mathbf{y})}{\partial x_n} \\ \frac{\partial f_2(\mathbf{x}, \mathbf{y})}{\partial x_1} & \frac{\partial f_2(\mathbf{x}, \mathbf{y})}{\partial x_2} & \cdots & \frac{\partial f_2(\mathbf{x}, \mathbf{y})}{\partial x_n} \\ \vdots & & \ddots & \\ \frac{\partial f_q(\mathbf{x}, \mathbf{y})}{\partial x_1} & \frac{\partial f_q(\mathbf{x}, \mathbf{y})}{\partial x_2} & \cdots & \frac{\partial f_q(\mathbf{x}, \mathbf{y})}{\partial x_n} \end{pmatrix}$$

while

$$D_2\mathbf{f}(\mathbf{x}, \mathbf{y}) = \begin{pmatrix} \frac{\partial f_1(\mathbf{x}, \mathbf{y})}{\partial y_1} & \frac{\partial f_1(\mathbf{x}, \mathbf{y})}{\partial y_2} & \cdots & \frac{\partial f_1(\mathbf{x}, \mathbf{y})}{\partial y_p} \\ \frac{\partial f_2(\mathbf{x}, \mathbf{y})}{\partial y_1} & \frac{\partial f_2(\mathbf{x}, \mathbf{y})}{\partial y_2} & \cdots & \frac{\partial f_2(\mathbf{x}, \mathbf{y})}{\partial y_p} \\ \vdots & & \ddots & \\ \frac{\partial f_q(\mathbf{x}, \mathbf{y})}{\partial y_1} & \frac{\partial f_q(\mathbf{x}, \mathbf{y})}{\partial y_2} & \cdots & \frac{\partial f_q(\mathbf{x}, \mathbf{y})}{\partial y_p} \end{pmatrix}.$$

There have been quite a few terms defined. First there was the concept of continuity. Next the concept of partial or directional derivative. Next there was the concept of differentiability and the derivative being a linear transformation determined by a certain matrix. Finally, it was shown that if a function is  $C^1$ , then it has a derivative. To give a rough idea of the relationships of these topics, here is a picture.



You might ask whether there are examples of functions which are differentiable but not  $C^1$ . Of course there are. In fact, such an example exists even in one dimension. Consider

$$f(x) = \begin{cases} x^2 \sin\left(\frac{1}{x}\right) & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases}. \quad (20.26)$$

Then you should verify that  $f'(x)$  exists for all  $x \in \mathbb{R}$  but  $f'$  fails to be continuous at  $x = 0$ . Thus the function is differentiable at every point of  $\mathbb{R}$  but fails to be  $C^1$  at every point of  $\mathbb{R}$ .

## 20.10 Exercises

1. In more advanced treatments of differentiation the following is done. Let  $\mathbf{f} : U \rightarrow \mathbb{R}^q$  where  $U$  is an open set in  $\mathbb{R}^p$ .  $\mathbf{f}$  is differentiable at  $\mathbf{x} \in U$  if there exists a linear transformation,  $L$  such that

$$\mathbf{f}(\mathbf{x} + \mathbf{v}) = \mathbf{f}(\mathbf{x}) + L\mathbf{v} + \mathbf{o}(\mathbf{v}).$$

Show that defining  $L \equiv D\mathbf{f}(\mathbf{x})$ , then  $D\mathbf{f}(\mathbf{x})$  is well defined. That is, there can't be two different linear transformations which satisfy the above equation. **Hint:** Replace  $\mathbf{v}$  with  $t\mathbf{v}$  and use the above to write

$$\frac{\mathbf{f}(\mathbf{x} + t\mathbf{v}) - \mathbf{f}(\mathbf{x})}{t} = L_i\mathbf{v} + \frac{\mathbf{o}(t)}{t}, i = 1, 2.$$

Now this implies  $L_1\mathbf{v} = L_2\mathbf{v} + \frac{\mathbf{o}(t)}{t}$ .

2. Suppose  $\mathbf{f} : U \rightarrow \mathbb{R}^q$  and let  $\mathbf{x} \in U$  and  $\mathbf{v}$  be a unit vector. Show  $D_{\mathbf{v}}\mathbf{f}(\mathbf{x}) = D\mathbf{f}(\mathbf{x})\mathbf{v}$ .
3. Now, using the definition in Problem 1, show that  $J\mathbf{f}(\mathbf{x})$  the above matrix of partial derivatives exists and  $D\mathbf{f}(\mathbf{x})\mathbf{v} = L\mathbf{v} = J\mathbf{f}(\mathbf{x})\mathbf{v}$ .
4. Let

$$f(x, y) = \begin{cases} x & \text{if } |y| > |x| \\ -x & \text{if } |y| \leq |x| \end{cases}.$$

Show  $f$  is continuous at  $(0, 0)$  and that the partial derivatives exist at  $(0, 0)$  but the function is not differentiable at  $(0, 0)$ .

5. Let

$$f(x, y) = \begin{cases} \frac{(x^2 - y^4)^2}{(x^2 + y^4)^2} & \text{if } (x, y) \neq (0, 0) \\ 1 & \text{if } (x, y) = (0, 0) \end{cases}.$$

Show that all directional derivatives of  $f$  exist at  $(0, 0)$ , and are all equal to zero but the function is not even continuous at  $(0, 0)$ .

## 20.11 The Gradient

Let  $f : U \rightarrow \mathbb{R}$  where  $U$  is an open subset of  $\mathbb{R}^n$  and suppose  $f$  is differentiable on  $U$ . Thus if  $\mathbf{x} \in U$ ,

$$f(\mathbf{x} + \mathbf{v}) = f(\mathbf{x}) + \sum_{j=1}^n \frac{\partial f(\mathbf{x})}{\partial x_j} v_j + o(\mathbf{v}). \quad (20.27)$$

**Definition 20.11.1** Define  $\nabla f(\mathbf{x}) \equiv \left( \frac{\partial f}{\partial x_1}(\mathbf{x}), \dots, \frac{\partial f}{\partial x_n}(\mathbf{x}) \right)^T$ . This vector is called the gradient vector.

This defines the gradient for a scalar valued function. There are ways to define the gradient for vector valued functions but this will not be attempted in this book.

It follows immediately from (20.27) that

$$f(\mathbf{x} + \mathbf{v}) = f(\mathbf{x}) + \nabla f(\mathbf{x}) \cdot \mathbf{v} + o(\mathbf{v}) \quad (20.28)$$

An important aspect of the gradient is its relation with the directional derivative. From (20.28), for  $\mathbf{v}$  a unit vector,

$$\begin{aligned} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} &= \nabla f(\mathbf{x}) \cdot \mathbf{v} + \frac{o(t\mathbf{v})}{t} \\ &= \nabla f(\mathbf{x}) \cdot \mathbf{v} + \frac{o(t)}{t}. \end{aligned}$$

Therefore, taking  $t \rightarrow 0$ ,

$$D_{\mathbf{v}}f(\mathbf{x}) = \nabla f(\mathbf{x}) \cdot \mathbf{v}. \quad (20.29)$$

**Example 20.11.2** Let  $f(x, y, z) = x^2 + \sin(xy) + z$ . Find  $D_{\mathbf{v}}f(1, 0, 1)$  where  $\mathbf{v} = \left( \frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}} \right)$ .

Note this vector which is given is already a unit vector. Therefore, from the above, it is only necessary to find  $\nabla f(1, 0, 1)$  and take the dot product.  $\nabla f(x, y, z) = (2x, x \cos(xy), 1)$ . Therefore,  $\nabla f(1, 0, 1) = (2, 1, 1)$ . Therefore, the directional derivative is  $(2, 1, 1) \cdot \left( \frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}} \right) = \frac{4}{3}\sqrt{3}$ .

Because of (20.29) it is easy to find the largest possible directional derivative and the smallest possible directional derivative.

**Proposition 20.11.3** Let  $f : U \rightarrow \mathbb{R}$  be a differentiable function and let  $\mathbf{x} \in U$ . Then

$$\max \{D_{\mathbf{v}}f(\mathbf{x}) : |\mathbf{v}| = 1\} = |\nabla f(\mathbf{x})| \quad (20.30)$$

and

$$\min \{D_{\mathbf{v}}f(\mathbf{x}) : |\mathbf{v}| = 1\} = -|\nabla f(\mathbf{x})|. \quad (20.31)$$

Furthermore, the maximum in (20.30) occurs when  $\mathbf{v} = \nabla f(\mathbf{x}) / |\nabla f(\mathbf{x})|$  and the minimum in (20.31) occurs when  $\mathbf{v} = -\nabla f(\mathbf{x}) / |\nabla f(\mathbf{x})|$ .

**Proof:** From (20.29) and the Cauchy Schwarz inequality,

$$|D_{\mathbf{v}}f(\mathbf{x})| \leq |\nabla f(\mathbf{x})|$$

and so for any choice of  $\mathbf{v}$  with  $|\mathbf{v}| = 1$ ,

$$-|\nabla f(\mathbf{x})| \leq D_{\mathbf{v}}f(\mathbf{x}) \leq |\nabla f(\mathbf{x})|.$$

The proposition is proved by noting that if  $\mathbf{v} = -\nabla f(\mathbf{x}) / |\nabla f(\mathbf{x})|$ , then

$$\begin{aligned} D_{\mathbf{v}}f(\mathbf{x}) &= \nabla f(\mathbf{x}) \cdot (-\nabla f(\mathbf{x}) / |\nabla f(\mathbf{x})|) \\ &= -|\nabla f(\mathbf{x})|^2 / |\nabla f(\mathbf{x})| = -|\nabla f(\mathbf{x})| \end{aligned}$$

while if  $\mathbf{v} = \nabla f(\mathbf{x}) / |\nabla f(\mathbf{x})|$ , then

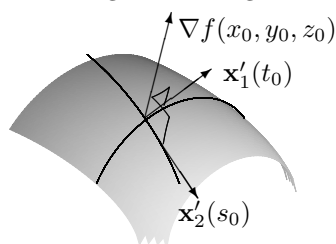
$$\begin{aligned} D_{\mathbf{v}}f(\mathbf{x}) &= \nabla f(\mathbf{x}) \cdot (\nabla f(\mathbf{x}) / |\nabla f(\mathbf{x})|) \\ &= |\nabla f(\mathbf{x})|^2 / |\nabla f(\mathbf{x})| = |\nabla f(\mathbf{x})|. \end{aligned}$$

The conclusion of the above proposition is important in many physical models. For example, consider some material which is at various temperatures depending on location. Because it has cool places and hot places, it is expected that the heat will flow from the hot places to the cool places. Consider a small surface having a unit normal,  $\mathbf{n}$ . Thus  $\mathbf{n}$  is a normal to this surface and has unit length. If it is desired to find the rate in calories per second at which heat crosses this little surface in the direction of  $\mathbf{n}$  it is defined as  $\mathbf{J} \cdot \mathbf{n}A$  where  $A$  is the area of the surface and  $\mathbf{J}$  is called the heat flux. It is reasonable to suppose the rate at which heat flows across this surface will be largest when  $\mathbf{n}$  is in the direction of greatest rate of decrease of the temperature. In other words, heat flows most readily in the direction which involves the maximum rate of decrease in temperature. This expectation will be realized by taking  $\mathbf{J} = -K\nabla u$  where  $K$  is a positive scalar function which can depend on a variety of things. The above relation between the heat flux and  $\nabla u$  is usually called the Fourier heat conduction law and the constant,  $K$  is known as the coefficient of thermal conductivity. It is a material property, different for iron than for aluminum. In most applications,  $K$  is considered to be a constant but this is wrong. Experiments show this scalar should depend on temperature. Nevertheless, things get very difficult if this dependence is allowed. The constant can depend on position in the material or even on time.

An identical relationship is usually postulated for the flow of a diffusing species. In this problem, something like a pollutant diffuses. It may be an insecticide in ground water for example. Like heat, it tries to move from areas of high concentration toward areas of low concentration. In this case  $\mathbf{J} = -K\nabla c$  where  $c$  is the concentration of the diffusing species. When applied to diffusion, this relationship is known as Fick's law. Mathematically, it is indistinguishable from the problem of heat flow.

Note the importance of the gradient in formulating these models.

The gradient has fundamental geometric significance illustrated by the following picture.



In this picture, the surface is a piece of a level surface of a function of three variables,  $f(x, y, z)$ . Thus the surface is defined by  $f(x, y, z) = c$  or more completely as  $\{(x, y, z) : f(x, y, z) = c\}$ . For example, if  $f(x, y, z) = x^2 + y^2 + z^2$ , this would be a piece of a sphere. There are two smooth curves in this picture which lie in the surface having parameterizations,  $\mathbf{x}_1(t) = (x_1(t), y_1(t), z_1(t))$  and  $\mathbf{x}_2(s) = (x_2(s), y_2(s), z_2(s))$  which intersect at the point,  $(x_0, y_0, z_0)$  on this surface<sup>4</sup>. This intersection occurs when  $t = t_0$  and  $s = s_0$ . Since the points,  $\mathbf{x}_1(t)$  for  $t$  in an interval lie in the level surface, it follows

<sup>4</sup>Do there exist any smooth curves which lie in the level surface of  $f$  and pass through the point  $(x_0, y_0, z_0)$ ? It turns out there do if  $\nabla f(x_0, y_0, z_0) \neq \mathbf{0}$  and if the function,  $f$ , is  $C^1$ . However, this is a

$$f(x_1(t), y_1(t), z_1(t)) = c$$

for all  $t$  in some interval. Therefore, taking the derivative of both sides and using the chain rule on the left,

$$\begin{aligned} & \frac{\partial f}{\partial x}(x_1(t), y_1(t), z_1(t)) x'_1(t) + \\ & \frac{\partial f}{\partial y}(x_1(t), y_1(t), z_1(t)) y'_1(t) + \frac{\partial f}{\partial z}(x_1(t), y_1(t), z_1(t)) z'_1(t) = 0. \end{aligned}$$

In terms of the gradient, this merely states

$$\nabla f(x_1(t), y_1(t), z_1(t)) \cdot \mathbf{x}'_1(t) = 0.$$

Similarly,

$$\nabla f(x_2(s), y_2(s), z_2(s)) \cdot \mathbf{x}'_2(s) = 0.$$

Letting  $s = s_0$  and  $t = t_0$ , it follows

$$\nabla f(x_0, y_0, z_0) \cdot \mathbf{x}'_1(t_0) = 0, \quad \nabla f(x_0, y_0, z_0) \cdot \mathbf{x}'_2(s_0) = 0.$$

It follows  $\nabla f(x_0, y_0, z_0)$  is perpendicular to both the direction vectors of the two indicated curves shown. Surely if things are as they should be, these two direction vectors would determine a plane which deserves to be called the tangent plane to the level surface of  $f$  at the point  $(x_0, y_0, z_0)$  and that  $\nabla f(x_0, y_0, z_0)$  is perpendicular to this tangent plane at the point,  $(x_0, y_0, z_0)$ .

**Example 20.11.4** Find the equation of the tangent plane to the level surface,  $f(x, y, z) = 6$  of the function,  $f(x, y, z) = x^2 + 2y^2 + 3z^2$  at the point  $(1, 1, 1)$ .

First note that  $(1, 1, 1)$  is a point on this level surface. Therefore, it suffices to find the normal vector to the proposed plane. But  $\nabla f(x, y, z) = (2x, 4y, 6z)$  and so  $\nabla f(1, 1, 1) = (2, 4, 6)$ . Therefore, from this problem, the equation of the plane is

$$(2, 4, 6) \cdot (x - 1, y - 1, z - 1) = 0$$

or in other words,

$$2x - 12 + 4y + 6z = 0.$$

## 20.12 Exercises

1. Find the gradients of  $f =$

- (a)  $x^2y + z^3$  at  $(1, 1, 2)$
- (b)  $z \sin(x^2y) + 2^{x+y}$  at  $(1, 1, 0)$
- (c)  $u \ln(x + y + z^2 + w)$  at  $(x, y, z, w, u) = (1, 1, 1, 1, 2)$

2. Find the directional derivatives of  $f$  at the indicated point in the direction,  $\left(\frac{1}{2}, \frac{1}{2}, \frac{1}{\sqrt{2}}\right)$ .

- (a)  $x^2y + z^3$  at  $(1, 1, 1)$

---

consequence of the implicit function theorem, one of the greatest theorems in all mathematics and a topic for an advanced calculus class.

- (b)  $z \sin(x^2 y) + 2^{x+y}$  at  $(1, 1, 2)$   
 (c)  $xy + z^2 + w$  at  $(1, 2, 3)$
3. Find the tangent plane to the indicated level surface at the indicated point.
- (a)  $x^2 y + z^3 = 2$  at  $(1, 1, 1)$   
 (b)  $z \sin(x^2 y) + 2^{x+y} = 2 \sin 1 + 4$  at  $(1, 1, 2)$   
 (c)  $\cos(x) + z \sin(x + y) = 1$  at  $(-\pi, \frac{3\pi}{2}, 2)$
4. Explain why the displacement vector of an object moving in  $\mathbb{R}^3$  is always perpendicular to the velocity vector if the object is always at a fixed distance from a given point.
5. The level surfaces  $x^2 + y^2 + z^2 = 4$  and  $z + x^2 + y^2 = 4$  have the point  $(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2}, 1)$  in the curve formed by the intersection of these surfaces. Find a direction vector for this curve at this point. **Hint:** Recall the gradients of the two surfaces are perpendicular to the corresponding surfaces at this point. A direction vector for the desired curve should be perpendicular to both of these gradients.
6. In a slightly more general setting, suppose  $f_1(x, y, z) = 0$  and  $f_2(x, y, z) = 0$  are two level surfaces which intersect in a curve which has parameterization,  $(x(t), y(t), z(t))$ . Find a differential equation for this curve.

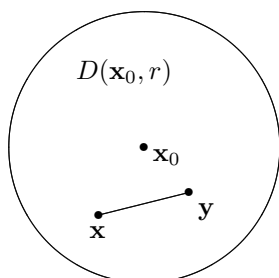
## 20.13 Nonlinear Ordinary Differential Equations Local Existence

Recall Definition 20.7.3. The following lemma applies to such functions.

**Lemma 20.13.1** *Let  $D(\mathbf{x}_0, r) \equiv \{\mathbf{x} \in \mathbb{R}^n : |\mathbf{x} - \mathbf{x}_0| \leq r\}$  and suppose  $U$  is an open set containing  $D(\mathbf{x}_0, r)$  such that  $\mathbf{f} : U \rightarrow \mathbb{R}^n$  is  $C^1(U)$ . Then for  $K = Mn$ , where  $M$  denotes the maximum of  $\left| \frac{\partial \mathbf{f}}{\partial x_i}(\mathbf{z}) \right|$  for  $\mathbf{z} \in D(\mathbf{x}_0, r)$  and  $i = 1, \dots, n$ , it follows that for all  $\mathbf{x}, \mathbf{y} \in D(\mathbf{x}_0, r)$ ,*

$$|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| \leq K |\mathbf{x} - \mathbf{y}|.$$

**Proof:** Let  $\mathbf{x}, \mathbf{y} \in D(\mathbf{x}_0, r)$  and consider the line segment joining these two points,  $\mathbf{x} + t(\mathbf{y} - \mathbf{x})$  for  $t \in [0, 1]$  as shown in the following picture.



Observe  $\mathbf{x} + (0)(\mathbf{y} - \mathbf{x}) = \mathbf{x}$  while  $\mathbf{x} + (1)(\mathbf{y} - \mathbf{x}) = \mathbf{y}$  and for  $t \in (0, 1)$ ,

$$\begin{aligned} |\mathbf{x} + t(\mathbf{y} - \mathbf{x}) - \mathbf{x}_0| &= |(1-t)(\mathbf{x} - \mathbf{x}_0) + t(\mathbf{y} - \mathbf{x}_0)| \\ &\leq (1-t)|\mathbf{x} - \mathbf{x}_0| + t|\mathbf{y} - \mathbf{x}_0| \\ &\leq (1-t)r + tr = r \end{aligned}$$



showing that the points on the line segment  $\mathbf{x} + t(\mathbf{y} - \mathbf{x})$  for  $t \in [0, 1]$  are contained in  $D(\mathbf{x}_0, r)$  as claimed and as indicated in the picture.

If  $\mathbf{h}(t) = \mathbf{f}(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))$  for  $t \in [0, 1]$ , then

$$\mathbf{f}(\mathbf{y}) - \mathbf{f}(\mathbf{x}) = \mathbf{h}(1) - \mathbf{h}(0) = \int_0^1 \mathbf{h}'(t) dt.$$

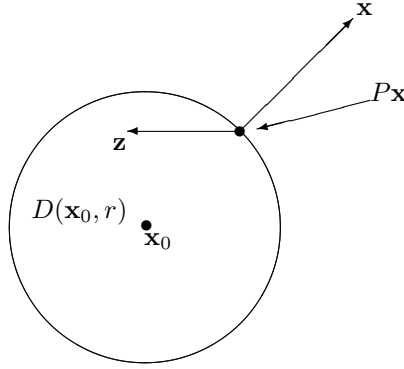
Also, by the chain rule,

$$\mathbf{h}'(t) = \sum_{i=1}^n \frac{\partial \mathbf{f}}{\partial x_i}(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))(y_i - x_i).$$

Therefore,

$$\begin{aligned} |\mathbf{f}(\mathbf{y}) - \mathbf{f}(\mathbf{x})| &= \left| \int_0^1 \sum_{i=1}^n \frac{\partial \mathbf{f}}{\partial x_i}(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))(y_i - x_i) dt \right| \\ &\leq \int_0^1 \sum_{i=1}^n \left| \frac{\partial \mathbf{f}}{\partial x_i}(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) \right| |y_i - x_i| dt \\ &\leq M \sum_{i=1}^n |y_i - x_i| \leq Mn |\mathbf{x} - \mathbf{y}|. \end{aligned}$$

Now consider the map,  $P$  which maps all of  $\mathbb{R}^n$  to  $D(\mathbf{x}_0, r)$  given as follows. For  $\mathbf{x} \in D(\mathbf{x}_0, r)$ ,  $P\mathbf{x} = \mathbf{x}$ . For  $\mathbf{x} \notin D(\mathbf{x}_0, r)$   $P\mathbf{x}$  will be the closest point in  $D(\mathbf{x}_0, r)$  to  $\mathbf{x}$ .



**Lemma 20.13.2** For any pair of points,  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ ,  $|P\mathbf{x} - P\mathbf{y}| \leq |\mathbf{x} - \mathbf{y}|$ .

**Proof:** From the above picture and  $\mathbf{z} \in D(\mathbf{x}_0, r)$  arbitrary, the angle between the vectors  $\mathbf{x} - P\mathbf{x}$  and  $\mathbf{z} - P\mathbf{x}$  is always greater than  $\pi/2$  radians. Therefore, the cosine of this angle is always negative. It follows that

$$(\mathbf{y} - P\mathbf{y}) \cdot (P\mathbf{x} - P\mathbf{y}) \leq 0$$

and

$$(\mathbf{x} - P\mathbf{x}) \cdot (P\mathbf{y} - P\mathbf{x}) \leq 0.$$

Thus  $(\mathbf{x} - P\mathbf{x}) \cdot (P\mathbf{x} - P\mathbf{y}) \geq 0$  and so if subtracting yields

$$(\mathbf{x} - P\mathbf{x} - (\mathbf{y} - P\mathbf{y})) \cdot (P\mathbf{x} - P\mathbf{y}) \geq 0$$

which implies

$$\begin{aligned} (\mathbf{x} - \mathbf{y}) \cdot (P\mathbf{x} - P\mathbf{y}) &\geq (P\mathbf{x} - P\mathbf{y}) \cdot (P\mathbf{x} - P\mathbf{y}) \\ &= |P\mathbf{x} - P\mathbf{y}|^2. \end{aligned}$$

Now apply the Cauchy Schwarz inequality to the left side of the above inequality to obtain

$$|\mathbf{x} - \mathbf{y}| |P\mathbf{x} - P\mathbf{y}| \geq |P\mathbf{x} - P\mathbf{y}|^2$$

which yields the claim of the lemma.

This implies the following local existence and uniqueness theorem.

**Theorem 20.13.3** *Let  $[a, b]$  be a closed interval and let  $U$  be an open subset of  $\mathbb{R}^n$ . Let  $\mathbf{f} : [a, b] \times U \rightarrow \mathbb{R}^n$  be continuous and suppose that for each  $t \in [a, b]$ , the map  $\mathbf{x} \rightarrow \frac{\partial \mathbf{f}}{\partial x_i}(t, \mathbf{x})$  is continuous. Also let  $\mathbf{x}_0 \in U$  and  $c \in [a, b]$ . Then there exists an interval,  $I \subseteq [a, b]$  such that  $c \in I$  and there exists a unique solution to the initial value problem,*

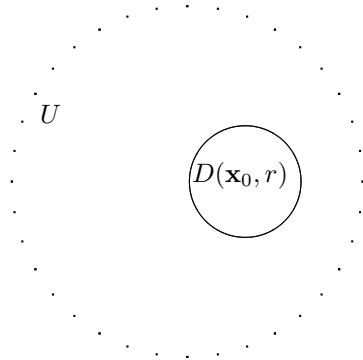
$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}), \quad \mathbf{x}(c) = \mathbf{x}_0 \quad (20.32)$$

valid for  $t \in I$ . If  $D(\mathbf{x}_0, r)$  is a closed ball as defined above which is contained in  $U$ , and if

$$M \equiv \max \{ |\mathbf{f}(t, \mathbf{x})| : (t, \mathbf{x}) \in [a, b] \times D(\mathbf{x}_0, r) \},$$

then one can take  $I = [p, q]$  where  $p = \max(a, c - \frac{r}{M})$  and  $q = \min(b, c + \frac{r}{M})$ .

**Proof:** Consider the following picture.



The large dotted circle represents  $U$  and the little solid circle represents  $D(\mathbf{x}_0, r)$  as indicated. Here  $r$  is chosen so small that  $D(\mathbf{x}_0, r)$  is contained in  $U$  as shown. Now let  $P$  denote the projection map defined above. Consider the initial value problem

$$\mathbf{x}' = \mathbf{f}(t, P\mathbf{x}), \quad \mathbf{x}(c) = \mathbf{x}_0. \quad (20.33)$$

From Lemma 20.13.1 and the continuity of  $\mathbf{x} \rightarrow \frac{\partial \mathbf{f}}{\partial x_i}(t, \mathbf{x})$ , there exists a constant,  $K$  such that if  $\mathbf{x}, \mathbf{y} \in D(\mathbf{x}_0, r)$ , then  $|\mathbf{f}(t, \mathbf{x}) - \mathbf{f}(t, \mathbf{y})| \leq K |\mathbf{x} - \mathbf{y}|$  for all  $t \in [a, b]$ . Therefore, by Lemma 20.13.2,

$$|\mathbf{f}(t, P\mathbf{x}) - \mathbf{f}(t, P\mathbf{y})| \leq K |P\mathbf{x} - P\mathbf{y}| \leq K |\mathbf{x} - \mathbf{y}|.$$

It follows from Theorem 15.7.6 on Page 370 that (20.33) has a unique solution valid for  $t \in [a, b]$ . Since  $\mathbf{x}$  is continuous, it follows that there exists an interval,  $I$  containing  $c$  such that for  $t \in I$ ,  $\mathbf{x}(t) \in B(\mathbf{x}_0, r)$ . Therefore, for these values of  $t$ ,  $\mathbf{f}(t, P\mathbf{x}) = \mathbf{f}(t, \mathbf{x})$  and so there is a unique solution to (20.32) on  $I$ .

It remains to describe  $I$ .

$$|\mathbf{x}(t) - \mathbf{x}_0| = \left| \int_c^t \mathbf{x}'(s) ds \right| \leq \left| \int_c^t |\mathbf{x}'(s)| ds \right| \leq M |t - c|$$

and so, if  $|t - c| < \frac{r}{M}$ , it follows  $\mathbf{x}(t) \in D(\mathbf{x}_0, r)$  and so for  $[p, q]$  as given, if  $t \in [p, q]$ ,  $P\mathbf{x}(t) = \mathbf{x}(t)$  and the solution to (20.33) is the solution to (20.32). This proves the theorem.

There is a more general theorem which just gives existence.

**Theorem 20.13.4** *Let  $[a, b]$  be a closed interval and let  $U$  be an open subset of  $\mathbb{R}^n$ . Let  $\mathbf{f} : [a, b] \times U \rightarrow \mathbb{R}^n$  be continuous. Also let  $\mathbf{x}_0 \in U$  and  $c \in [a, b]$ . Then there exists an interval,  $I \subseteq [a, b]$  such that  $c \in I$  and there exists a solution to the initial value problem,*

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}), \quad \mathbf{x}(c) = \mathbf{x}_0 \quad (20.34)$$

*valid for  $t \in I$ . If  $D(\mathbf{x}_0, r)$  is a closed ball which is contained in  $U$ , and if*

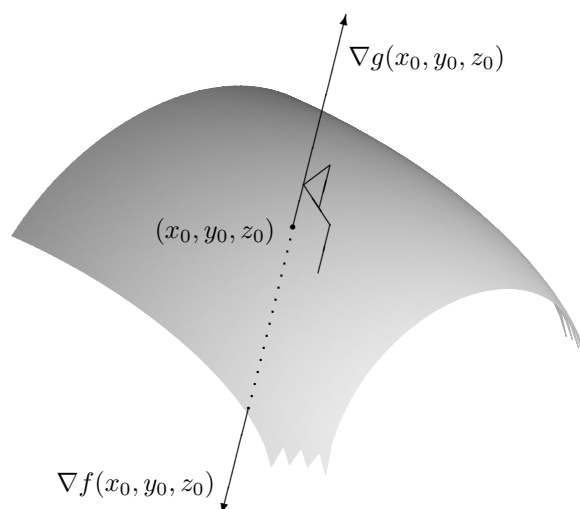
$$M \equiv \max \{ |\mathbf{f}(t, \mathbf{x})| : (t, \mathbf{x}) \in [a, b] \times D(\mathbf{x}_0, r) \},$$

*then one can take  $I = [p, q]$  where  $p = \max(a, c - \frac{r}{M})$  and  $q = \min(b, c + \frac{r}{M})$ .*

This is called the Peano (pronounced piano) existence theorem and it is one of the most useful results in certain areas of applied math. The proof is not particularly hard and may be found as an exercise in [10] or [14]. However, it does involve the notion of compactness in a space of continuous functions and this is a topic for an advanced calculus course.

## 20.14 Lagrange Multipliers

Lagrange multipliers are used to solve extremum problems for a function defined on a level set of another function. For example, suppose you want to maximize  $xy$  given that  $x + y = 4$ . This is not too hard to do using methods developed earlier. Solve for one of the variables, say  $y$ , in the constraint equation,  $x + y = 4$  to find  $y = 4 - x$ . Then the function to maximize is  $f(x) = x(4 - x)$  and the answer is clearly  $x = 2$ . Thus the two numbers are  $x = y = 2$ . This was easy because you could easily solve the constraint equation for one of the variables in terms of the other. Now what if you wanted to maximize  $f(x, y, z) = xyz$  subject to the constraint that  $x^2 + y^2 + z^2 = 4$ ? It is still possible to do this using the techniques of Problem 5 on Page 473. Solve for one of the variables in the constraint equation, say  $z$ , substitute it into  $f$ , and then find where the partial derivatives equal zero to find candidates for the extremum. However, it seems you might encounter many cases and it does look a little fussy. However, sometimes you can't solve the constraint equation for one variable in terms of the others. Also, what if you had many constraints. What if you wanted to maximize  $f(x, y, z)$  subject to the constraints  $x^2 + y^2 = 4$  and  $z = 2x + 3y^2$ . Things are clearly getting more involved and messy. It turns out that at an extremum, there is a simple relationship between the gradient of the function to be maximized and the gradient of the constraint function. This relation can be seen geometrically as in the following picture.



In the picture, the surface represents a piece of the level surface of  $g(x, y, z) = 0$  and  $f(x, y, z)$  is the function of three variables which is being maximized or minimized on the level surface and suppose the extremum of  $f$  occurs at the point  $(x_0, y_0, z_0)$ . As shown above,  $\nabla g(x_0, y_0, z_0)$  is perpendicular to the surface or more precisely to the tangent plane. However, if  $\mathbf{x}(t) = (x(t), y(t), z(t))$  is a point on a smooth curve which passes through  $(x_0, y_0, z_0)$  when  $t = t_0$ , then the function,  $h(t) = f(x(t), y(t), z(t))$  must have either a maximum or a minimum at the point,  $t = t_0$ . Therefore,  $h'(t_0) = 0$ . But this means

$$\begin{aligned} 0 &= h'(t_0) = \nabla f(x(t_0), y(t_0), z(t_0)) \cdot \mathbf{x}'(t_0) \\ &= \nabla f(x_0, y_0, z_0) \cdot \mathbf{x}'(t_0) \end{aligned}$$

and since this holds for any such smooth curve,  $\nabla f(x_0, y_0, z_0)$  is also perpendicular to the surface. This picture represents a situation in three dimensions and you can see that it is intuitively clear that this implies  $\nabla f(x_0, y_0, z_0)$  is some scalar multiple of  $\nabla g(x_0, y_0, z_0)$ . Thus

$$\nabla f(x_0, y_0, z_0) = \lambda \nabla g(x_0, y_0, z_0)$$

This  $\lambda$  is called a Lagrange multiplier after Lagrange who considered such problems in the 1700's.

Of course the above argument is at best only heuristic. It does not deal with the question of existence of smooth curves lying in the constraint surface passing through  $(x_0, y_0, z_0)$ . Nor does it consider all cases, being essentially confined to three dimensions. In addition to this, it fails to consider the situation in which there are many constraints. However, I think it is likely a geometric notion like that presented above which led Lagrange to formulate the method.

**Example 20.14.1** Maximize  $xyz$  subject to  $x^2 + y^2 + z^2 = 27$ .

Here  $f(x, y, z) = xyz$  while  $g(x, y, z) = x^2 + y^2 + z^2 - 27$ . Then  $\nabla g(x, y, z) = (2x, 2y, 2z)$  and  $\nabla f(x, y, z) = (yz, xz, xy)$ . Then at the point which maximizes this function<sup>5</sup>,

$$(yz, xz, xy) = \lambda (2x, 2y, 2z).$$

<sup>5</sup>There exists such a point because the sphere is closed and bounded.

Therefore, each of  $2\lambda x^2, 2\lambda y^2, 2\lambda z^2$  equals  $xyz$ . It follows that at any point which maximizes  $xyz$ ,  $|x| = |y| = |z|$ . Therefore, the only candidates for the point where the maximum occurs are  $(3, 3, 3), (-3, -3, 3), (-3, 3, 3)$ , etc. The maximum occurs at  $(3, 3, 3)$  which can be verified by plugging in to the function which is being maximized.

There are no magic bullets here. It was still required to solve a system of nonlinear equations to get the answer. However, it does often help to do it this way.

The above generalizes to a general procedure. First it is essential to have a special case of the implicit function theorem.

**Theorem 20.14.2** *Let  $\mathbf{f} : (a, b) \times U \rightarrow \mathbb{R}^n$  where  $U$  is an open set in  $\mathbb{R}^n$  and suppose  $D_1\mathbf{f}$  and  $D_2\mathbf{f}$  are both continuous on  $(a, b) \times U$ . Suppose also that for  $t_0 \in (a, b)$  and  $\mathbf{x}_0 \in U$ ,*

$$\mathbf{f}(t_0, \mathbf{x}_0) = \mathbf{0}, \det(D_2\mathbf{f}(t_0, \mathbf{x}_0)) \neq 0.$$

*Then there exists  $(p, q)$  such that  $t_0 \in (p, q)$  and a function,  $\mathbf{x} : (p, q) \rightarrow U$  such that*

$$\mathbf{f}(t, \mathbf{x}(t)) = \mathbf{0}$$

*for all  $t \in (p, q)$ .*

**Proof:** By the Peano existence theorem<sup>6</sup>, Theorem 20.13.4 on Page 499 there exists an open interval,  $(p, q)$  containing  $t_0$  and a solution,  $\mathbf{x}$  to the initial value problem,

$$\mathbf{x}'(t) = -D_2\mathbf{f}(t, \mathbf{x}(t))^{-1} D_1\mathbf{f}(t, \mathbf{x}(t)), \mathbf{x}(t_0) = \mathbf{x}_0$$

for  $t \in (p, q)$ . This is because by assumption, the function,  $(t, \mathbf{x}) \rightarrow -D_2\mathbf{f}(t, \mathbf{x})^{-1} D_1\mathbf{f}(t, \mathbf{x})$  is continuous. Thus,

$$D_1\mathbf{f}(t, \mathbf{x}(t)) + D_2\mathbf{f}(t, \mathbf{x}(t)) \mathbf{x}'(t) = \mathbf{0}.$$

Then

$$\mathbf{0} = \mathbf{f}(t_0, \mathbf{x}_0) \tag{20.35}$$

and by the chain rule,

$$\frac{d}{dt}(\mathbf{f}(t, \mathbf{x}(t))) = D_1\mathbf{f}(t, \mathbf{x}(t)) + D_2\mathbf{f}(t, \mathbf{x}(t)) \mathbf{x}'(t) = \mathbf{0}$$

Therefore,  $t \rightarrow \mathbf{f}(t, \mathbf{x}(t))$  must be a constant vector and this vector equals zero by (20.35). This proves the theorem.

Let  $f : U \rightarrow \mathbb{R}$  be a  $C^1$  function and let

$$g_i(\mathbf{x}) = 0, \quad i = 1, \dots, m \tag{20.36}$$

be a collection of equality constraints with  $m < n$ . Now consider the system of nonlinear equations

$$\begin{aligned} f(\mathbf{x}) &= t \\ g_i(\mathbf{x}) &= 0, \quad i = 1, \dots, m. \end{aligned}$$

$\mathbf{x}_0$  is a local maximum if  $f(\mathbf{x}_0) \geq f(\mathbf{x})$  for all  $\mathbf{x}$  near  $\mathbf{x}_0$  which also satisfies the constraints (20.36). A local minimum is defined similarly.

<sup>6</sup>I have not presented a proof of this important theorem so this part of the argument is technically incomplete. However you can just assume in addition that the function  $(t, \mathbf{x}) \rightarrow -D_2\mathbf{f}(t, \mathbf{x}(t))^{-1} D_1\mathbf{f}(t, \mathbf{x}(t))$  is  $C^1$  and then use the Theorem on the local existence for solutions to ordinary differential equations which was proved. Not much is lost in doing this because this extra regularity will be available in nearly every conceivable example of any interest.

Now suppose  $\mathbf{x}_0$  satisfies all the constraints and is a point where  $f$  has a local maximum. (The case of a local minimum is handled similarly.) Thus  $t_0 \equiv f(\mathbf{x}_0)$  is the largest value of any  $f(\mathbf{x})$  for  $\mathbf{x}$  near  $\mathbf{x}_0$  satisfying the constraints. Consider the  $m+1 \times n$  matrix,

$$\begin{pmatrix} f_{x_1}(\mathbf{x}_0) & f_{x_2}(\mathbf{x}_0) & \cdots & f_{x_n}(\mathbf{x}_0) \\ g_{1x_1}(\mathbf{x}_0) & g_{1x_2}(\mathbf{x}_0) & \cdots & g_{1x_n}(\mathbf{x}_0) \\ \vdots & \vdots & & \vdots \\ g_{mx_1}(\mathbf{x}_0) & g_{mx_2}(\mathbf{x}_0) & \cdots & g_{mx_n}(\mathbf{x}_0) \end{pmatrix}. \quad (20.37)$$

It will be shown that this matrix has rank less than  $m+1$ . This will be established by deriving a contradiction if the matrix has rank  $m+1$ .

Suppose therefore, this matrix has rank  $m+1$  then by Definition 19.9.19 on Page 456, some  $m+1 \times m+1$  submatrix has nonzero determinant. This determines  $m+1$  variables,  $x_{i_1}, \dots, x_{i_{m+1}}$  corresponding to the columns used to get the  $m+1 \times m+1$  matrix having nonzero determinant. For  $\mathbf{x} \equiv (x_1, \dots, x_n)$ , if  $j$  is not equal to one of the  $i_k$ , replace the variable,  $x_j$  with the constant  $x_{0j}$  where  $\mathbf{x}_0 \equiv (x_{01}, \dots, x_{0n})$  and let  $\mathbf{y} \in \mathbb{R}^{m+1}$  be defined as

$$\mathbf{y} \equiv (y_1, \dots, y_{m+1}) \equiv (x_{i_1}, \dots, x_{i_{m+1}})$$

and let

$$\mathbf{y}_0 \equiv (y_{01}, \dots, y_{0m+1}) \equiv (x_{0i_1}, \dots, x_{0i_{m+1}}).$$

Denote by  $W$  the set of points in  $\mathbb{R}^{m+1}$

$$W \equiv \{\mathbf{y} \in \mathbb{R}^{m+1} : \mathbf{y} = (x_{i_1}, \dots, x_{i_{m+1}}) \text{ where } \mathbf{x} \in U.\}$$

Now let  $f^1(\mathbf{y})$  be defined as what is obtained from  $f(\mathbf{x})$  by replacing  $x_j$  for  $j$  none of the  $i_k$  with  $x_{0j}$  and  $x_{i_k}$  with  $y_k$  and let  $g^l(\mathbf{y})$  be defined similarly for  $l = 1, 2, \dots, m$ . Therefore,

$$f^1(\mathbf{y}_0) = t_0, \quad g^l(\mathbf{y}_0) = 0,$$

and for all  $\mathbf{y}$  close enough to  $\mathbf{y}_0$  with  $g^l(\mathbf{y}) = 0$  for each  $l = 1, \dots, m$ ,  $t_0 \geq f^1(\mathbf{y})$ . Furthermore,

$$\begin{pmatrix} f_{y_1}^1(\mathbf{y}_0) & f_{y_2}^1(\mathbf{y}_0) & \cdots & f_{y_{m+1}}^1(\mathbf{y}_0) \\ g_{y_1}^1(\mathbf{y}_0) & g_{y_2}^1(\mathbf{y}_0) & \cdots & g_{y_{m+1}}^1(\mathbf{y}_0) \\ \vdots & \vdots & & \vdots \\ g_{y_1}^m(\mathbf{y}_0) & g_{y_2}^m(\mathbf{y}_0) & \cdots & g_{y_{m+1}}^m(\mathbf{y}_0) \end{pmatrix}$$

has nonzero determinant. But this matrix equals  $D_2\mathbf{F}(t_0, \mathbf{y}_0)$  where  $\mathbf{F} : \mathbb{R} \times W \rightarrow \mathbb{R}^{m+1}$  is defined as

$$\mathbf{F}(t, \mathbf{y}) \equiv \begin{pmatrix} f^1(\mathbf{y}) - t \\ g^1(\mathbf{y}) \\ \vdots \\ g^m(\mathbf{y}) \end{pmatrix} \quad (20.38)$$

Therefore by Theorem 20.14.2 there is an open interval containing  $t_0$ ,  $(p, q)$ , and a function  $\mathbf{y}$  defined on this interval with the property that  $\mathbf{F}(t, \mathbf{y}(t)) = 0$  for all  $t \in (p, q)$ . This is a contradiction because, as noted above,  $t_0$  is the largest value for all  $f^1(\mathbf{y})$  with  $\mathbf{y}$  close to  $\mathbf{y}_0$  such that  $g^l(\mathbf{y}) = 0$ . Therefore, the matrix of (20.37) has rank less than  $m+1$  as claimed.

Suppose now the rank of the matrix,

$$\begin{pmatrix} g_{1x_1}(\mathbf{x}_0) & g_{1x_2}(\mathbf{x}_0) & \cdots & g_{1x_n}(\mathbf{x}_0) \\ \vdots & \vdots & & \vdots \\ g_{mx_1}(\mathbf{x}_0) & g_{mx_2}(\mathbf{x}_0) & \cdots & g_{mx_n}(\mathbf{x}_0) \end{pmatrix} \quad (20.39)$$

equals  $m$ . This means there is an  $m \times m$  submatrix which has nonzero determinant and so none of the above rows can be a linear combination of the others. Therefore, by Theorem 19.9.20 on Page 456 every row of the matrix of (20.37) is a linear combination of the rows of the above matrix. In particular, there exist scalars,

$$\lambda_1, \dots, \lambda_m$$

such that

$$\begin{pmatrix} f_{x_1}(\mathbf{x}_0) \\ \vdots \\ f_{x_n}(\mathbf{x}_0) \end{pmatrix} = \lambda_1 \begin{pmatrix} g_{1x_1}(\mathbf{x}_0) \\ \vdots \\ g_{1x_n}(\mathbf{x}_0) \end{pmatrix} + \dots + \lambda_m \begin{pmatrix} g_{mx_1}(\mathbf{x}_0) \\ \vdots \\ g_{mx_n}(\mathbf{x}_0) \end{pmatrix} \quad (20.40)$$

The  $\lambda_i$  are called the Lagrange multipliers. This proves the following theorem.

**Theorem 20.14.3** *Let  $U$  be an open subset of  $\mathbb{R}^n$  and let  $f : U \rightarrow \mathbb{R}$  be a  $C^1$  function. Then if  $\mathbf{x}_0 \in U$  is either a local maximum or local minimum of  $f$  subject to the constraints (20.36), and if the rank of the matrix in (20.39) equals  $m$ , then there exist scalars,  $\lambda_1, \dots, \lambda_m$  such that (20.40) holds.*

To help remember how to use (20.40) it may be helpful to do the following. First write the Lagrangian,

$$L = f(\mathbf{x}) - \sum_{i=1}^m \lambda_i g_i(\mathbf{x})$$

and then proceed to take derivatives with respect to each of the components of  $\mathbf{x}$  and also derivatives with respect to each  $\lambda_i$  and set all of these equations equal to 0. The formula (20.40) is what results from taking the derivatives of  $L$  with respect to the components of  $\mathbf{x}$ . When you take the derivatives with respect to the Lagrange multipliers, and set what results equal to 0, you just pick up the constraint equations. This yields  $n + m$  equations for the  $n + m$  unknowns,  $x_1, \dots, x_n, \lambda_1, \dots, \lambda_m$ . Then you proceed to look for solutions to these equations. Of course these might be impossible to find using methods of algebra, but you just do your best and hope it will work out.

**Example 20.14.4** *Minimize  $xyz$  subject to the constraints  $x^2 + y^2 + z^2 = 4$  and  $x - 2y = 0$ .*

Form the Lagrangian,

$$L = xyz - \lambda(x^2 + y^2 + z^2 - 4) - \mu(x - 2y)$$

and proceed to take derivatives with respect to every possible variable, leading to the following system of equations.

$$\begin{aligned} yz - 2\lambda x - \mu &= 0 \\ xz - 2\lambda y + 2\mu &= 0 \\ xy - 2\lambda z &= 0 \\ x^2 + y^2 + z^2 &= 4 \\ x - 2y &= 0 \end{aligned}$$

Now you have to find the solutions to this system of equations. In general, this could be very hard or even impossible. If  $\lambda = 0$ , then from the third equation, either  $x$  or  $y$  must equal 0. Therefore, from the first two equations,  $\mu = 0$  also. If  $\mu = 0$  and  $\lambda \neq 0$ , then from the first two equations,  $xyz = 2\lambda x^2$  and  $xyz = 2\lambda y^2$  and so either  $x = y$  or  $x = -y$ , which

requires that both  $x$  and  $y$  equal zero thanks to the last equation. But then from the fourth equation,  $z = \pm 2$  and now this contradicts the third equation. Thus  $\mu$  and  $\lambda$  are either both equal to zero or neither one is and the expression,  $xyz$  equals zero in this case. However, I know this is not the best value for a minimizer because I can take  $x = 2\sqrt{\frac{3}{5}}, y = \sqrt{\frac{3}{5}}$ , and  $z = -1$ . This satisfies the constraints and the product of these numbers equals a negative number. Therefore, both  $\mu$  and  $\lambda$  must be non zero. Now use the last equation eliminate  $x$  and write the following system.

$$\begin{aligned} 5y^2 + z^2 &= 4 \\ y^2 - \lambda z &= 0 \\ yz - \lambda y + \mu &= 0 \\ yz - 4\lambda y - \mu &= 0 \end{aligned}$$

From the last equation,  $\mu = (yz - 4\lambda y)$ . Substitute this into the third and get

$$\begin{aligned} 5y^2 + z^2 &= 4 \\ y^2 - \lambda z &= 0 \\ yz - \lambda y + yz - 4\lambda y &= 0 \end{aligned}$$

$y = 0$  will not yield the minimum value from the above example. Therefore, divide the last equation by  $y$  and solve for  $\lambda$  to get  $\lambda = (2/5)z$ . Now put this in the second equation to conclude

$$\begin{aligned} 5y^2 + z^2 &= 4 \\ y^2 - (2/5)z^2 &= 0 \end{aligned}$$

a system which is easy to solve. Thus  $y^2 = 8/15$  and  $z^2 = 4/3$ . Therefore, candidates for minima are  $(2\sqrt{\frac{8}{15}}, \sqrt{\frac{8}{15}}, \pm\sqrt{\frac{4}{3}})$ , and  $(-2\sqrt{\frac{8}{15}}, -\sqrt{\frac{8}{15}}, \pm\sqrt{\frac{4}{3}})$ , a choice of 4 points to check.

Clearly the one which gives the smallest value is  $(2\sqrt{\frac{8}{15}}, \sqrt{\frac{8}{15}}, -\sqrt{\frac{4}{3}})$  or  $(-2\sqrt{\frac{8}{15}}, -\sqrt{\frac{8}{15}}, -\sqrt{\frac{4}{3}})$  and the minimum value of the function subject to the constraints is  $-\frac{2}{5}\sqrt{30} - \frac{2}{3}\sqrt{3}$ .

You should rework this problem first solving the second easy constraint for  $x$  and then producing a simpler problem involving only the variables  $y$  and  $z$ .

## 20.15 Exercises

1. Maximize  $2x + 3y - 6z$  subject to the constraint,  $x^2 + 2y^2 + 3z^2 = 9$ .
2. Find the dimensions of the largest rectangle which can be inscribed in a circle of radius  $r$ .
3. Find the points on  $y^2x = 9$  which are closest to  $(0, 0)$ .
4. Find points on  $xy = 4$  farthest from  $(0, 0)$  if any exist. If none exist, tell why. What does this say about the method of Lagrange multipliers?
5. A can is supposed to have a volume of  $36\pi$  cubic centimeters. Find the dimensions of the can which minimizes the surface area.
6. A can is supposed to have a volume of  $36\pi$  cubic centimeters. The top and bottom of the can are made of tin costing 4 cents per square centimeter and the sides of the can are made of aluminum costing 5 cents per square centimeter. Find the dimensions of the can which minimizes the cost.



7. Minimize  $\sum_{j=1}^n x_j$  subject to the constraint  $\sum_{j=1}^n x_j^2 = a^2$ . Your answer should be some function of  $a$  which you may assume is a positive number.
8. Find the point,  $(x, y, z)$  on the level surface,  $4x^2 + y^2 - z^2 = 1$  which is closest to  $(0, 0, 0)$ .
9. A curve is formed from the intersection of the plane,  $2x + 3y + z = 3$  and the cylinder  $x^2 + y^2 = 4$ . Find the point on this curve which is closest to  $(0, 0, 0)$ .
10. A curve is formed from the intersection of the plane,  $2x + 3y + z = 3$  and the sphere  $x^2 + y^2 + z^2 = 16$ . Find the point on this curve which is closest to  $(0, 0, 0)$ .
11. Find the point on the plane,  $2x + 3y + z = 4$  which is closest to the point  $(1, 2, 3)$ .
12. Let  $A = (A_{ij})$  be an  $n \times n$  matrix which is symmetric. Thus  $A_{ij} = A_{ji}$  and recall  $(A\mathbf{x})_i = A_{ij}x_j$  where as usual sum over the repeated index. Show  $\frac{\partial}{\partial x_i}(A_{ij}x_jx_i) = 2A_{ij}x_j$ . Show that when you use the method of Lagrange multipliers to maximize the function,  $A_{ij}x_jx_i$  subject to the constraint,  $\sum_{j=1}^n x_j^2 = 1$ , the value of  $\lambda$  which corresponds to the maximum value of this functions is such that  $A_{ij}x_j = \lambda x_i$ . Thus  $A\mathbf{x} = \lambda\mathbf{x}$ . This  $\lambda$  is called an eigenvalue of the matrix,  $A$ .
13. Here are two lines.  $\mathbf{x} = (1 + 2t, 2 + t, 3 + t)^T$  and  $\mathbf{x} = (2 + s, 1 + 2s, 1 + 3s)^T$ . Find points  $\mathbf{p}_1$  on the first line and  $\mathbf{p}_2$  on the second with the property that  $|\mathbf{p}_1 - \mathbf{p}_2|$  is at least as small as the distance between any other pair of points, one chosen on one line and the other on the other line.
14. Find the dimensions of the largest triangle which can be inscribed in a circle of radius  $r$ .
15. Find the point on the intersection of  $z = x^2 + y^2$  and  $x + y + z = 1$  which is closest to  $(0, 0, 0)$ .
16. Minimize  $4x^2 + y^2 + 9z^2$  subject to  $x + y - z = 1$  and  $x - 2y + z = 0$ .
17. Minimize  $xyz$  subject to the constraints  $x^2 + y^2 + z^2 = r^2$  and  $x - y = 0$ .
18. Let  $n$  be a positive integer. Find  $n$  numbers whose sum is  $8n$  and the sum of the squares is as small as possible.
19. Find the point on the level surface,  $2x^2 + xy + z^2 = 16$  which is closest to  $(0, 0, 0)$ .
20. Find the point on  $\frac{x^2}{4} + \frac{y^2}{9} + z^2 = 1$  closest to the plane  $x + y + z = 10$ .
21. Let  $x_1, \dots, x_5$  be 5 positive numbers. Maximize their product subject to the constraint that

$$x_1 + 2x_2 + 3x_3 + 4x_4 + 5x_5 = 300.$$

22. Let  $f(x_1, \dots, x_n) = x_1^n x_2^{n-1} \cdots x_n^1$ . Then  $f$  achieves a maximum on the set,

$$S \equiv \left\{ \mathbf{x} \in \mathbb{R}^n : \sum_{i=1}^n ix_i = 1 \text{ and each } x_i \geq 0 \right\}.$$

If  $\mathbf{x} \in S$  is the point where this maximum is achieved, find  $x_1/x_n$ .

23. Let  $(x, y)$  be a point on the ellipse,  $x^2/a^2 + y^2/b^2 = 1$  which is in the first quadrant. Extend the tangent line through  $(x, y)$  till it intersects the  $x$  and  $y$  axes and let  $A(x, y)$  denote the area of the triangle formed by this line and the two coordinate axes. Find the maximum value of the area of this triangle as a function of  $a$  and  $b$ .
24. Maximize  $\prod_{i=1}^n x_i^2$  ( $\equiv x_1^2 \times x_2^2 \times x_3^2 \times \cdots \times x_n^2$ ) subject to the constraint,  $\sum_{i=1}^n x_i^2 = r^2$ . Show the maximum is  $(r^2/n)^n$ . Now show from this that

$$\left( \prod_{i=1}^n x_i^2 \right)^{1/n} \leq \frac{1}{n} \sum_{i=1}^n x_i^2$$

and finally, conclude that if each number  $x_i \geq 0$ , then

$$\left( \prod_{i=1}^n x_i \right)^{1/n} \leq \frac{1}{n} \sum_{i=1}^n x_i$$

and there exist values of the  $x_i$  for which equality holds. This says the “geometric mean” is always smaller than the arithmetic mean.

25. Maximize  $x^2 y^2$  subject to the constraint

$$\frac{x^{2p}}{p} + \frac{y^{2q}}{q} = r^2$$

where  $p, q$  are real numbers larger than 1 which have the property that

$$\frac{1}{p} + \frac{1}{q} = 1.$$

show the maximum is achieved when  $x^{2p} = y^{2q}$  and equals  $r^2$ . Now conclude that if  $x, y > 0$ , then

$$xy \leq \frac{x^p}{p} + \frac{y^q}{q}$$

and there are values of  $x$  and  $y$  where this inequality is an equation.

## 20.16 Taylor's Formula For Functions Of Many Variables

Let  $U$  be an open subset of  $\mathbb{R}^p$  and let  $f : U \rightarrow \mathbb{R}$  be a  $C^k$  function. What is the analog of Taylor's formula for a function of one variable? To answer this question, let  $\mathbf{x} \in U$  and let  $B(\mathbf{x}, 2r) \subseteq U$ . Then if  $\mathbf{v} \in \mathbb{R}^p$  with  $|\mathbf{v}| < r$ , it follows that  $\mathbf{x} + t\mathbf{v} \in U$  whenever  $|t| \leq 2$ . The consideration of this question involves the function

$$F(t) \equiv f(\mathbf{x} + t\mathbf{v})$$

for  $t \in (0, 2)$ . This is a function of one variable and so the one variable Taylor's formula applies to  $F$  provided  $F$  has derivatives. That this is the case is part of the next lemma. First a definition is convenient.

**Definition 20.16.1** Let  $\partial_i$  denote the partial derivative with respect to  $x_i$  where  $\mathbf{x} = (x_1, \dots, x_p)$ . Thus

$$\partial_i g \equiv \frac{\partial g}{\partial x_i}.$$

Now for  $\mathbf{v} = (v_1, \dots, v_p)$ , a fixed vector, define the "differential operator",  $\sum_{i=1}^p v_i \partial_i$  according to the rule

$$\begin{aligned} \left( \sum_{i=1}^p v_i \partial_i \right) (g) &\equiv \sum_{i=1}^p v_i (\partial_i g) \\ &= \sum_{i=1}^p v_i \frac{\partial g}{\partial x_i}. \end{aligned}$$

The symbol,  $(\sum_{i=1}^p v_i \partial_i)^2$  is defined just as you might expect.

$$\left( \sum_{i=1}^p v_i \partial_i \right)^2 = \sum_{i,j} v_i v_j \partial_i \partial_j$$

where

$$\left( \sum_{i,j} v_i v_j \partial_i \partial_j \right) (g) = \sum_{i,j} v_i v_j (\partial_i (\partial_j g)).$$

Similarly,

$$\left( \sum_{i=1}^p v_i \partial_i \right)^3 = \sum_{i,j,k} v_i v_j v_k \partial_i \partial_j \partial_k$$

and more generally,

$$\left( \sum_{i=1}^p v_i \partial_i \right)^n = \sum_{i_1, \dots, i_n} v_{i_1} v_{i_2} \cdots v_{i_n} \partial_{i_1} \partial_{i_2} \cdots \partial_{i_n}$$

where the sum is taken over all choices of the indices,  $i_1, \dots, i_n$  in  $\{1, \dots, p\}$ .

**Lemma 20.16.2** Let  $f$  and  $F$  be as above. Then  $F \in C^k(0, 2)$  and for each  $l \leq k$

$$F^{(l)}(t) = \left( \sum_{i=1}^p v_i \partial_i \right)^l f(\mathbf{x} + t\mathbf{v}).$$

**Proof:** From the chain rule,

$$F'(t) = \sum_{i=1}^p v_i \partial_i f(\mathbf{x} + t\mathbf{v}) = \left( \sum_{i=1}^p v_i \partial_i \right)^1 f(\mathbf{x} + t\mathbf{v})$$

and this proves the lemma if  $l = 1$ . Now suppose this is true for  $l < k$ . Then

$$\begin{aligned} F^{(l)}(t) &= \left( \sum_{i=1}^p v_i \partial_i \right)^l f(\mathbf{x} + t\mathbf{v}) \\ &= \sum_{i_1, \dots, i_l} v_{i_1} v_{i_2} \cdots v_{i_l} \partial_{i_1} \partial_{i_2} \cdots \partial_{i_l} f(\mathbf{x} + t\mathbf{v}) \end{aligned}$$

Taking the derivative of this using the chain rule, yields

$$\begin{aligned} F^{(l+1)}(t) &= \sum_{i_1, \dots, i_l} v_{i_1} v_{i_2} \cdots v_{i_l} \partial_{i_1} \partial_{i_2} \cdots \partial_{i_l} \left( \sum_{i_{l+1}=1}^p v_{i_{l+1}} \partial_{i_{l+1}} \right) f(\mathbf{x} + t\mathbf{v}) \\ &= \sum_{i_1, \dots, i_{l+1}} v_{i_1} v_{i_2} \cdots v_{i_{l+1}} \partial_{i_1} \partial_{i_2} \cdots \partial_{i_{l+1}} f(\mathbf{x} + t\mathbf{v}) \\ &= \left( \sum_{i=1}^p v_i \partial_i \right)^{l+1} f(\mathbf{x} + t\mathbf{v}) \end{aligned}$$

and this proves the lemma.

By Taylor's theorem, Theorem 11.1.1 on Page 257,

$$F(t) = F(0) + F'(0)t + \cdots + \frac{F^{(k-1)}(0)t^{k-1}}{(k-1)!} + \frac{F^{(k)}(s_t)t^k}{k!}$$

for some  $s_t \in (0, t)$ . In particular, this formula holds for  $t = 1$  since  $1 < 2$ . Therefore, from the above lemma, there exists  $s \in (0, 1)$  such that

$$\begin{aligned} f(\mathbf{x} + \mathbf{v}) &= f(\mathbf{x}) + \left( \sum_{i=1}^p v_i \partial_i \right) f(\mathbf{x}) + \cdots \\ &\quad \cdots + \frac{1}{(k-1)!} \left( \sum_{i=1}^p v_i \partial_i \right)^{k-1} f(\mathbf{x}) + \frac{1}{k!} \left( \sum_{i=1}^p v_i \partial_i \right)^k f(\mathbf{x} + s\mathbf{v}) \end{aligned} \quad (20.41)$$

and this is Taylor's formula for a function of many variables. This proves the following theorem.

**Theorem 20.16.3** *Let  $U$  be an open set in  $\mathbb{R}^p$  and let  $f : U \rightarrow \mathbb{R}$  be  $C^k$ . Then there exists  $s \in (0, 1)$  such that (20.41) holds whenever  $|\mathbf{v}|$  is sufficiently small.*

### 20.16.1 Some Linear Algebra

If  $H$  is an  $n \times n$  matrix with  $H = H^T$ , it follows from Exercise 8 on Page 420

$$H\mathbf{x} \cdot \mathbf{y} = \mathbf{x} \cdot H\mathbf{y} = H\mathbf{y} \cdot \mathbf{x}. \quad (20.42)$$

In case you didn't work this exercise, here is a proof of this exercise.

**Lemma 20.16.4** *Let  $A$  be an  $m \times n$  matrix. Then  $A\mathbf{x} \cdot \mathbf{y} = \mathbf{x} \cdot A^T \mathbf{y}$ .*

**Proof:** Using the repeated index summation convention,

$$A\mathbf{x} \cdot \mathbf{y} = \mathbf{A}_{ij} x_j y_i = x_j A_{ij} y_i = x_j (A^T)_{ji} y_i = \mathbf{x} \cdot A^T \mathbf{y}.$$

This lemma implies (20.42).

**Theorem 20.16.5** *Let  $H$  be a real  $n \times n$  symmetric matrix and let  $S$  denote the vectors in  $\mathbb{R}^n$  which have unit length. Then there exists  $\lambda, \mu$  such that*

$$\lambda = \min \{ H\mathbf{x} \cdot \mathbf{x} : \mathbf{x} \in S \}, \quad \mu = \max \{ H\mathbf{x} \cdot \mathbf{x} : \mathbf{x} \in S \}.$$

*If  $\mathbf{x}_\lambda$  is a point of  $S$  at which  $H\mathbf{x}_\lambda \cdot \mathbf{x}_\lambda = \lambda$ , then  $H\mathbf{x}_\lambda = \lambda\mathbf{x}_\lambda$ . Similarly, if  $\mathbf{x}_\mu$  is a point of  $S$  at which  $H\mathbf{x}_\mu \cdot \mathbf{x}_\mu = \mu$ , then  $H\mathbf{x}_\mu = \mu\mathbf{x}_\mu$ . In addition to this, if  $H\mathbf{x} = r\mathbf{x}$  for some real number,  $r$  and  $\mathbf{x} \neq \mathbf{0}$ , then  $\lambda \leq r \leq \mu$ .*

**Proof:** A vector in  $S$  determines a unique point,  $\mathbf{x}$  which is at unit distance from  $\mathbf{0}$ . The set of such points,  $\mathbf{x}$  is a closed and bounded set in  $\mathbb{R}^n$  and so since  $\mathbf{x} \rightarrow H\mathbf{x} \cdot \mathbf{x}$  is continuous, it achieves a minimum value,  $\lambda$  and a maximum value,  $\mu$ . It remains to show  $\mu$  and  $\lambda$  satisfy the given conditions. First note that if  $\mathbf{u}$  is any nonzero vector in  $\mathbb{R}^n$ ,

$$H \frac{\mathbf{u}}{|\mathbf{u}|} \cdot \frac{\mathbf{u}}{|\mathbf{u}|} \geq \lambda$$

and so

$$H\mathbf{u} \cdot \mathbf{u} \geq \lambda |\mathbf{u}|^2 = \lambda \mathbf{u} \cdot \mathbf{u}$$

and so  $(H\mathbf{u} - \lambda\mathbf{u}) \cdot \mathbf{u} \geq 0$  for all  $\mathbf{u} \neq \mathbf{0}$ .

Now let  $\mathbf{v} \in \mathbb{R}^n$  be an arbitrary vector. Then from what was just shown

$$H(\mathbf{x}_\lambda + t\mathbf{v}) \cdot (\mathbf{x}_\lambda + t\mathbf{v}) \geq \lambda(\mathbf{x}_\lambda + t\mathbf{v}) \cdot (\mathbf{x}_\lambda + t\mathbf{v})$$

and so from properties of dot products and (20.42),

$$H\mathbf{x}_\lambda \cdot \mathbf{x}_\lambda + 2tH\mathbf{x}_\lambda \cdot \mathbf{v} + t^2H\mathbf{v} \cdot \mathbf{v} \geq \overbrace{\lambda\mathbf{x}_\lambda \cdot \mathbf{x}_\lambda}^{=\lambda} + 2t\lambda\mathbf{x}_\lambda \cdot \mathbf{v} + t^2\lambda\mathbf{v} \cdot \mathbf{v}.$$

Since  $H\mathbf{x}_\lambda \cdot \mathbf{x}_\lambda = \lambda$ , subtract this from both sides and then divide both sides by  $t$ . If  $t$  is negative this yields

$$2H\mathbf{x}_\lambda \cdot \mathbf{v} + tH\mathbf{v} \cdot \mathbf{v} \leq 2\lambda\mathbf{x}_\lambda \cdot \mathbf{v} + t\lambda\mathbf{v} \cdot \mathbf{v}$$

Now let  $t \rightarrow 0$  to obtain

$$2H\mathbf{x}_\lambda \cdot \mathbf{v} \leq 2\lambda\mathbf{x}_\lambda \cdot \mathbf{v}.$$

Letting  $t > 0$  and doing the same thing leads to

$$2H\mathbf{x}_\lambda \cdot \mathbf{v} \geq 2\lambda\mathbf{x}_\lambda \cdot \mathbf{v}.$$

Consequently,  $H\mathbf{x}_\lambda \cdot \mathbf{v} = \lambda\mathbf{x}_\lambda \cdot \mathbf{v}$  and so  $(H - \lambda I)\mathbf{x}_\lambda \cdot \mathbf{v} = 0$ . Since this holds for any  $\mathbf{v}$ , it holds in particular for  $\mathbf{v} = (H - \lambda I)\mathbf{x}_\lambda$  and therefore,  $(H - \lambda I)\mathbf{x}_\lambda = 0$ . The claim about  $\mathbf{x}_\mu$  and  $\mu$  is completely similar.

Now suppose  $H\mathbf{x} = r\mathbf{x}$ . This means  $H \frac{\mathbf{x}}{|\mathbf{x}|} \cdot \frac{\mathbf{x}}{|\mathbf{x}|} = r \frac{\mathbf{x}}{|\mathbf{x}|} \cdot \frac{\mathbf{x}}{|\mathbf{x}|} = r$  and so  $\lambda \leq r \leq \mu$  by the way these two numbers are defined. This proves the theorem.

### 20.16.2 The Second Derivative Test

Let  $f : U \rightarrow \mathbb{R}$  where  $U$  is an open set in  $\mathbb{R}^n$  and  $f$  is  $C^2$ . What is the analog of the second derivative test for functions of one variable? Let  $\mathbf{x}$  be a point of  $U$  where  $Df(\mathbf{x}) = \mathbf{0}$ . Thus all partial derivatives of  $f$  at this point equal zero. These points are called critical points. By Theorem 20.16.3 the following equation is valid for all  $\mathbf{v}$  having  $|\mathbf{v}|$  small enough.

$$\begin{aligned} f(\mathbf{x} + \mathbf{v}) &= f(\mathbf{x}) + \frac{1}{2} \left( \sum_{i=1}^p v_i \partial_i \right)^2 f(\mathbf{x} + s\mathbf{v}) \\ &= f(\mathbf{x}) + \frac{1}{2} \sum_{i,j} v_i v_j \partial_i \partial_j f(\mathbf{x} + s\mathbf{v}) \\ &= f(\mathbf{x}) + \frac{1}{2} \mathbf{v}^T H(\mathbf{x} + s\mathbf{v}) \mathbf{v} \\ &= f(\mathbf{x}) + \frac{1}{2} H(\mathbf{x} + s\mathbf{v}) \mathbf{v} \cdot \mathbf{v} \end{aligned} \tag{20.43}$$

where

$$H(\mathbf{y}) \equiv \begin{pmatrix} f_{x_1 x_1}(\mathbf{y}) & f_{x_1 x_2}(\mathbf{y}) & \cdots & f_{x_1 x_n}(\mathbf{y}) \\ f_{x_2 x_1}(\mathbf{y}) & f_{x_2 x_2}(\mathbf{y}) & \cdots & f_{x_2 x_n}(\mathbf{y}) \\ \vdots & \vdots & \ddots & \vdots \\ f_{x_n x_1}(\mathbf{y}) & f_{x_n x_2}(\mathbf{y}) & \cdots & f_{x_n x_n}(\mathbf{y}) \end{pmatrix}$$

and this is a symmetric matrix due to Theorem 20.4.1 on Page 474. This matrix is called the Hessian matrix. By the assumption  $f$  is  $C^2$ , all the entries of  $H(\mathbf{y})$  are continuous functions of  $\mathbf{y}$ . The following theorem is the second derivative test for a function of many variables.

**Theorem 20.16.6** *Let  $f : U \rightarrow \mathbb{R}$  for  $U$  an open set in  $\mathbb{R}^n$  and let  $f$  be a  $C^2$  function and suppose that at some  $\mathbf{x} \in U$ ,  $\partial_i f(\mathbf{x}) = 0$  for all  $i$ . Also let  $\mu$  and  $\lambda$  be defined for  $H = H(\mathbf{x})$  as in Theorem 20.16.5. Thus*

$$\lambda = \min \{ H(\mathbf{x}) \mathbf{u} \cdot \mathbf{u} : \mathbf{u} \in S \}, \quad \mu = \max \{ H(\mathbf{x}) \mathbf{u} \cdot \mathbf{u} : \mathbf{u} \in S \}.$$

*If  $\lambda > 0$  then  $f$  has a local minimum at  $\mathbf{x}$ . If  $\mu < 0$  then  $f$  has a local maximum at  $\mathbf{x}$ . If either  $\lambda$  or  $\mu$  equals zero, the test fails. If  $\lambda < 0$  and  $\mu > 0$  there exists a direction in which when  $f$  is evaluated on the line through the critical point having this direction, the resulting function of one variable has a local minimum and there exists a direction in which when  $f$  is evaluated on the line through the critical point having this direction, the resulting function of one variable has a local maximum. This last case is called a saddle point.*

**Proof:** Suppose first  $\lambda > 0$ . From (20.43),

$$f(\mathbf{x} + \mathbf{v}) = f(\mathbf{x}) + \frac{1}{2} H(\mathbf{x}) \mathbf{v} \cdot \mathbf{v} + \frac{1}{2} (H(\mathbf{x} + s\mathbf{v}) - H(\mathbf{x})) \mathbf{v} \cdot \mathbf{v} \quad (20.44)$$

By continuity of the entries of  $H(\mathbf{y})$ , it follows that if  $|\mathbf{v}|$  is small enough,

$$|(H(\mathbf{x} + s\mathbf{v}) - H(\mathbf{x})) \mathbf{v} \cdot \mathbf{v}| \leq \frac{\lambda}{2} |\mathbf{v}|^2$$

Therefore, for all  $\mathbf{v}$  small enough,

$$\begin{aligned} f(\mathbf{x} + \mathbf{v}) &\geq f(\mathbf{x}) + \frac{1}{2} H(\mathbf{x}) \frac{\mathbf{v}}{|\mathbf{v}|} \cdot \frac{\mathbf{v}}{|\mathbf{v}|} |\mathbf{v}|^2 + \frac{1}{2} \left( -\frac{\lambda}{2} |\mathbf{v}|^2 \right) \\ &\geq f(\mathbf{x}) + \frac{\lambda}{4} |\mathbf{v}|^2 \end{aligned}$$

showing  $f$  has a local minimum at  $\mathbf{x}$ .

The case where  $\mu < 0$  is completely similar and is left as an exercise. It remains to verify the test fails in case either  $\mu$  or  $\lambda = 0$ . This is easily done by looking at  $f(x, y) = x^4 + y^4$  in which  $f$  clearly has a local minimum at  $(0, 0)$ ,  $f(x, y) = -x^4 - y^4$  in which  $f$  has a local maximum at  $(0, 0)$ , and  $f(x, y) = x^4 - y^4$  in which  $(0, 0)$  is clearly neither a local minimum nor a local maximum for this function.

Finally consider the case where  $\mu > 0$  and  $\lambda < 0$ . By Theorem 20.16.5 there exists  $\mathbf{v}$  such that  $H(\mathbf{x}) \mathbf{v} = \mu \mathbf{v}$ . Thus  $H(\mathbf{x}) \mathbf{v} \cdot \mathbf{v} = \mu |\mathbf{v}|^2$ . The line through  $\mathbf{x}$  having  $\mathbf{v}$  as its direction vector is just  $\mathbf{x} + t\mathbf{v}$  and so from (20.44)

$$f(\mathbf{x} + t\mathbf{v}) = f(\mathbf{x}) + \frac{t^2}{2} H(\mathbf{x}) \mathbf{v} \cdot \mathbf{v} + \frac{t^2}{2} (H(\mathbf{x} + s\mathbf{v}) - H(\mathbf{x})) \mathbf{v} \cdot \mathbf{v} \quad (20.45)$$

$$= f(\mathbf{x}) + \frac{\mu t^2}{2} |\mathbf{v}|^2 + \frac{t^2}{2} (H(\mathbf{x} + s\mathbf{v}) - H(\mathbf{x})) \mathbf{v} \cdot \mathbf{v}. \quad (20.46)$$

Now consider the last term. From continuity of the partial derivatives,

$$|(H(\mathbf{x} + s\mathbf{v}) - H(\mathbf{x})) \mathbf{v} \cdot \mathbf{v}| \leq \frac{\mu}{4} |\mathbf{v}|^2$$

provided  $|\mathbf{v}|$  is small enough. The line is unchanged if  $\mathbf{v}$  is multiplied by a small scalar. therefore, for small enough  $|\mathbf{v}|$  in a certain direction,

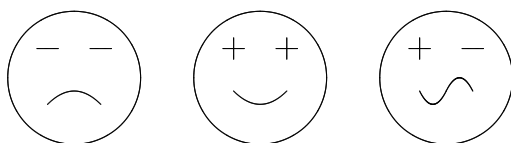
$$\begin{aligned} f(\mathbf{x} + t\mathbf{v}) &= f(\mathbf{x}) + \frac{\mu t^2}{2} |\mathbf{v}|^2 + \frac{t^2}{2} (H(\mathbf{x} + s\mathbf{v}) - H(\mathbf{x})) \mathbf{v} \cdot \mathbf{v} \\ &\geq f(\mathbf{x}) + \frac{\mu t^2}{2} |\mathbf{v}|^2 - \frac{\mu}{4} |\mathbf{v}|^2 = f(\mathbf{x}) + \frac{\mu t^2}{4} |\mathbf{v}|^2, \end{aligned}$$

showing that in that direction,  $f$  has a local minimum at the point,  $\mathbf{x}$ . The case where  $\lambda < 0$  is similar.

How can you tell which case occurs? The way to do this is to recall Theorem 20.16.5 in which it was shown that  $H\mathbf{x}_\lambda = \lambda\mathbf{x}_\lambda$  and  $H\mathbf{x}_\mu = \mu\mathbf{x}_\mu$ . In general, if  $H\mathbf{x} = \lambda\mathbf{x}$  and  $\mathbf{x} \neq \mathbf{0}$ , then  $(\lambda I - H)\mathbf{x} = \mathbf{0}$  even though  $\mathbf{x} \neq \mathbf{0}$ . Therefore, the matrix  $(\lambda I - H)$  must not have an inverse since if it did, you could multiply both sides of  $(\lambda I - H)\mathbf{x} = \mathbf{0}$  by this inverse and conclude  $\mathbf{x} = \mathbf{0}$ . Therefore,

$$\det(\lambda I - H) = 0.$$

It follows the numbers,  $\lambda$  and  $\mu$  will be found among the solutions to the above equation, called the characteristic equation of the matrix. The number  $\lambda$  is the smallest solution to this equation and the number  $\mu$  is the largest. The solutions to this equation are called eigenvalues of the matrix. Note this equation will always be an  $n^{\text{th}}$  degree polynomial equation. As in the case of a function of one variable, there is a picture to help you remember the situation.



In the picture the plus signs denote all positive eigenvalues. This is the case  $\lambda > 0$ . The minus signs denote all negative eigenvalues. This is the case  $\mu < 0$ .

**Example 20.16.7** Let  $f(x, y) = 2x^4 - 4x^3 + 14x^2 + 12yx^2 - 12yx - 12x + 2y^2 + 4y + 2$ . Find the critical points and determine whether they are local minimums, local maximums, or saddle points.

$f_x(x, y) = 8x^3 - 12x^2 + 28x + 24yx - 12y - 12$  and  $f_y(x, y) = 12x^2 - 12x + 4y + 4$ . The points at which both  $f_x$  and  $f_y$  equal zero are  $(\frac{1}{2}, -\frac{1}{4})$ ,  $(0, -1)$ , and  $(1, -1)$ .

The Hessian matrix is

$$\begin{pmatrix} 24x^2 + 28 + 24y - 24x & 24x - 12 \\ 24x - 12 & 4 \end{pmatrix}.$$

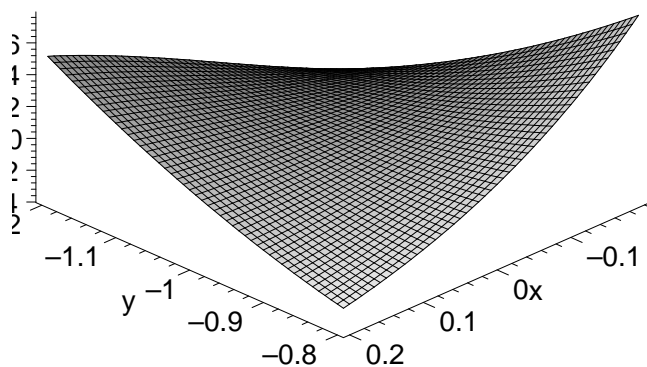
and the thing to determine is the sign of its eigenvalues evaluated at the critical points.

First consider the point  $(\frac{1}{2}, -\frac{1}{4})$ . This matrix is  $\begin{pmatrix} 16 & 0 \\ 0 & 4 \end{pmatrix}$  and its eigenvalues are 16, 4 showing that this is a local minimum.

Next consider  $(0, -1)$  at this point the Hessian matrix is  $\begin{pmatrix} 4 & -12 \\ -12 & 4 \end{pmatrix}$  and the eigenvalues are 16, -8. Therefore, this point is a saddle point.

Finally consider the point  $(1, -1)$ . At this point the Hessian is  $\begin{pmatrix} 4 & 12 \\ 12 & 4 \end{pmatrix}$  and the eigenvalues are  $16, -8$  so this point is also a saddle point.

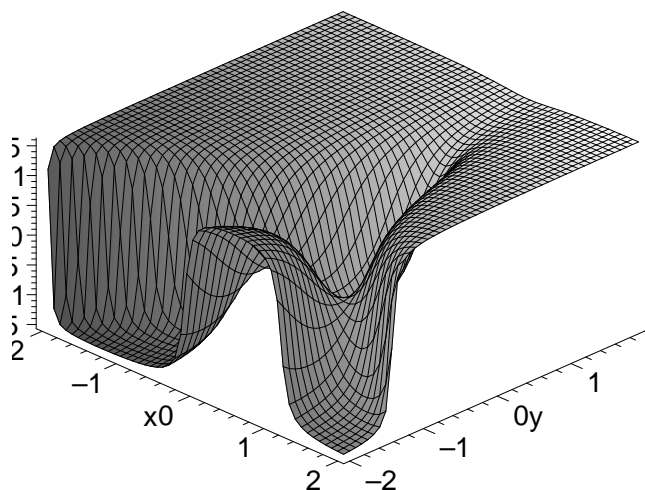
The geometric significance of a saddle point was explained above. In one direction it looks like a local minimum while in another it looks like a local maximum. In fact, they do look like a saddle. Here is a picture of the graph of the above function near the saddle point,  $(0, -1)$ .



You see it is a lot like the place where you sit on a saddle. If you want to get a better picture, you could graph instead

$$f(x, y) = \arctan(2x^4 - 4x^3 + 14x^2 + 12yx^2 - 12yx - 12x + 2y^2 + 4y + 2).$$

Since  $\arctan$  is a strictly increasing function, it preserves all the information about whether the given function is increasing or decreasing in certain directions. Below is a graph of this function which illustrates the behavior near the point  $(1, -1)$ .

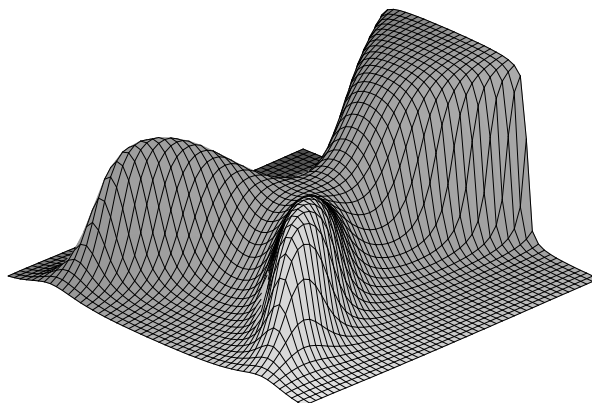


Of course sometimes the second derivative test is inadequate to determine what is going on. This should be no surprise since this was the case even for a function of one variable. For a function of two variables, a nice example is the Monkey saddle.

**Example 20.16.8** Suppose  $f(x, y) = \arctan(6xy^2 - 2x^3 - 3y^4)$ . Show  $(0, 0)$  is a critical point for which the second derivative test gives no information.

Before doing anything it might be interesting to look at the graph of this function of two variables plotted using Maple.





This picture should indicate why this is called a monkey saddle. It is because the monkey can sit in the saddle and have a place for his tail. Now to see  $(0, 0)$  is a critical point, note that

$$\frac{\partial (\arctan (g(x, y)))}{\partial x} = \frac{1}{1 + g(x, y)^2} g_x(x, y)$$

and that a similar formula holds for the partial derivative with respect to  $y$ . Therefore, it suffices to verify that for

$$g(x, y) = 6xy^2 - 2x^3 - 3y^4$$

$$g_x(0, 0) = g_y(0, 0) = 0.$$

$$g_x(x, y) = 6y^2 - 6x^2, \quad g_y(x, y) = 12xy - 12y^3$$

and clearly  $(0, 0)$  is a critical point. So are  $(1, 1)$  and  $(1, -1)$ . Now  $g_{xx}(0, 0) = 0$  and so does  $g_{xy}(0, 0)$  and  $g_{yy}(0, 0)$ . This implies  $f_{xx}, f_{xy}, f_{yy}$  are all equal to zero at  $(0, 0)$  also. (Why?) Therefore, the Hessian matrix is the zero matrix and clearly has only the zero eigenvalue. Therefore, the second derivative test is totally useless at this point.

However, suppose you took  $x = t$  and  $y = t$  and evaluated this function on this line. This reduces to  $h(t) = f(t, t) = \arctan(4t^3 - t^4)$ , which is strictly increasing near  $t = 0$ . This shows the critical point,  $(0, 0)$  of  $f$  is neither a local max. nor a local min. Next let  $x = 0$  and  $y = t$ . Then  $p(t) \equiv f(0, t) = -3t^4$ . Therefore, along the line,  $(0, t)$ ,  $f$  has a local maximum at  $(0, 0)$ .

## 20.17 Exercises

1. Finish the proof of Theorem 20.16.5.
2. Use the second derivative test on the critical points  $(1, 1)$ , and  $(1, -1)$  for Example 20.16.8.
3. If  $H = H^T$  and  $H\mathbf{x} = \lambda\mathbf{x}$  while  $H\mathbf{x} = \mu\mathbf{x}$  for  $\lambda \neq \mu$ , show  $\mathbf{x} \cdot \mathbf{y} = 0$ .
4. Verify the claims made about the three examples in Theorem 20.16.6.
5. Show the points  $(\frac{1}{2}, -\frac{21}{4})$ ,  $(0, -4)$ , and  $(1, -4)$  are critical points of the following function of three variables and classify them as local minimums, local maximums or saddle points.

$$f(x, y) = -x^4 + 2x^3 + 39x^2 + 10yx^2 - 10yx - 40x - y^2 - 8y - 16.$$

Answer:

The Hessian matrix is

$$\begin{pmatrix} -12x^2 + 78 + 20y + 12x & 20x - 10 \\ 20x - 10 & -2 \end{pmatrix}$$

The eigenvalues must be checked at the critical points. First consider the point  $(\frac{1}{2}, -\frac{21}{4})$ . At this point, the Hessian is

$$\begin{pmatrix} -24 & 0 \\ 0 & -2 \end{pmatrix}$$

and its eigenvalues are  $-24, -2$ , both negative. Therefore, the function has a local maximum at this point.

Next consider  $(0, -4)$ . At this point the Hessian matrix is

$$\begin{pmatrix} -2 & -10 \\ -10 & -2 \end{pmatrix}$$

and the eigenvalues are  $8, -12$  so the function has a saddle point.

Finally consider the point  $(1, -4)$ . The Hessian equals

$$\begin{pmatrix} -2 & 10 \\ 10 & -2 \end{pmatrix}$$

having eigenvalues:  $8, -12$  and so there is a saddle point here.

6. Show the points  $(\frac{1}{2}, -\frac{53}{12})$ ,  $(0, -4)$ , and  $(1, -4)$  are critical points of the following function of three variables and classify them according to whether they are local minimums, local maximums or saddle points.

$$f(x, y) = -3x^4 + 6x^3 + 37x^2 + 10yx^2 - 10yx - 40x - 3y^2 - 24y - 48.$$

Answer:

The Hessian matrix is

$$\begin{pmatrix} -36x^2 + 74 + 20y + 36x & 20x - 10 \\ 20x - 10 & -6 \end{pmatrix}.$$

Check its eigenvalues at the critical points. First consider the point  $(\frac{1}{2}, -\frac{53}{12})$ . At this point the Hessian is

$$\begin{pmatrix} -\frac{16}{3} & 0 \\ 0 & -6 \end{pmatrix}$$

and its eigenvalues are  $-\frac{16}{3}, -6$  so there is a local maximum at this point. The same analysis shows there are saddle points at the other two critical points.

7. Show the points  $(\frac{1}{2}, \frac{37}{20})$ ,  $(0, 2)$ , and  $(1, 2)$  are critical points of the following function of three variables and classify them according to whether they are local minimums, local maximums or saddle points.

$$f(x, y) = 5x^4 - 10x^3 + 17x^2 - 6yx^2 + 6yx - 12x + 5y^2 - 20y + 20.$$

Answer:

The Hessian matrix is

$$\begin{pmatrix} 60x^2 + 34 - 12y - 60x & -12x + 6 \\ -12x + 6 & 10 \end{pmatrix}.$$

Check its eigenvalues at the critical points. First consider the point  $(\frac{1}{2}, \frac{37}{20})$ . At this point, the Hessian matrix is

$$\begin{pmatrix} -\frac{16}{5} & 0 \\ 0 & 10 \end{pmatrix}$$

and its eigenvalues are  $-\frac{16}{5}, 10$ . Therefore, there is a saddle point.

Next consider  $(0, 2)$  at this point the Hessian matrix is

$$\begin{pmatrix} 10 & 6 \\ 6 & 10 \end{pmatrix}$$

and the eigenvalues are 16, 4. Therefore, there is a local minimum at this point. There is also a local minimum at the critical point,  $(1, 2)$ .

8. Show the points  $(\frac{1}{2}, -\frac{17}{8})$ ,  $(0, -2)$ , and  $(1, -2)$  are critical points of the following function of three variables and classify them according to whether they are local minimums, local maximums or saddle points.

$$f(x, y) = 4x^4 - 8x^3 - 4yx^2 + 4yx + 8x - 4x^2 + 4y^2 + 16y + 16.$$

Answer:

The Hessian matrix is  $\begin{pmatrix} 48x^2 - 8 - 8y - 48x & -8x + 4 \\ -8x + 4 & 8 \end{pmatrix}$ . Check its eigenvalues at the critical points. First consider the point  $(\frac{1}{2}, -\frac{17}{8})$ . This matrix is

$$\begin{pmatrix} -3 & 0 \\ 0 & 8 \end{pmatrix} \text{ and its eigenvalues are } -3, 8.$$

Next consider  $(0, -2)$  at this point the Hessian matrix is

$$\begin{pmatrix} 8 & 4 \\ 4 & 8 \end{pmatrix} \text{ and the eigenvalues are } 12, 4. \text{ Finally consider the point } (1, -2).$$

$$\begin{pmatrix} 8 & -4 \\ -4 & 8 \end{pmatrix}, \text{ eigenvalues: } 12, 4.$$

If the eigenvalues are both negative, then local max. If both positive, then local min. Otherwise the test fails.

9. Find the critical points of the following function of three variables and classify them according to whether they are local minimums, local maximums or saddle points.

$$f(x, y, z) = \frac{1}{3}x^2 + \frac{32}{3}x + \frac{4}{3} - \frac{16}{3}yx - \frac{58}{3}y - \frac{4}{3}zx - \frac{46}{3}z + \frac{1}{3}y^2 - \frac{4}{3}zy - \frac{5}{3}z^2.$$

Answer:

The critical point is at  $(-2, 3, -5)$ . The eigenvalues of the Hessian matrix at this point are  $-6, -2$ , and  $6$ .

10. Find the critical points of the following function of three variables and classify them according to whether they are local minimums, local maximums or saddle points.

$$f(x, y, z) = -\frac{5}{3}x^2 + \frac{2}{3}x - \frac{2}{3} + \frac{8}{3}yx + \frac{2}{3}y + \frac{14}{3}zx - \frac{28}{3}z - \frac{5}{3}y^2 + \frac{14}{3}zy - \frac{8}{3}z^2.$$

Answer:

The eigenvalues are 4,  $-10$ , and  $-6$  and the only critical point is  $(1, 1, 0)$ .

11. Find the critical points of the following function of three variables and classify them according to whether they are local minimums, local maximums or saddle points.

$$f(x, y, z) = -\frac{11}{3}x^2 + \frac{40}{3}x - \frac{56}{3} + \frac{8}{3}yx + \frac{10}{3}y - \frac{4}{3}zx + \frac{22}{3}z - \frac{11}{3}y^2 - \frac{4}{3}zy - \frac{5}{3}z^2.$$

12. Find the critical points of the following function of three variables and classify them according to whether they are local minimums, local maximums or saddle points.

$$f(x, y, z) = -\frac{2}{3}x^2 + \frac{28}{3}x + \frac{37}{3} + \frac{14}{3}yx + \frac{10}{3}y - \frac{4}{3}zx - \frac{26}{3}z - \frac{2}{3}y^2 - \frac{4}{3}zy + \frac{7}{3}z^2.$$

13. Show that if  $f$  has a critical point and some eigenvalue of the Hessian matrix is positive, then there exists a direction in which when  $f$  is evaluated on the line through the critical point having this direction, the resulting function of one variable has a local minimum. State and prove a similar result in the case where some eigenvalue of the Hessian matrix is negative.

14. Suppose  $\mu = 0$  but there are negative eigenvalues of the Hessian at a critical point. Show by giving examples that the second derivative tests fails.

15. Show the points  $(\frac{1}{2}, -\frac{9}{2})$ ,  $(0, -5)$ , and  $(1, -5)$  are critical points of the following function of three variables and classify them as local minimums, local maximums or saddle points.

$$f(x, y) = 2x^4 - 4x^3 + 42x^2 + 8yx^2 - 8yx - 40x + 2y^2 + 20y + 50.$$

16. Show the points  $(1, -\frac{11}{2})$ ,  $(0, -5)$ , and  $(2, -5)$  are critical points of the following function of three variables and classify them as local minimums, local maximums or saddle points.

$$f(x, y) = 4x^4 - 16x^3 - 4x^2 - 4yx^2 + 8yx + 40x + 4y^2 + 40y + 100.$$

17. Show the points  $(\frac{3}{2}, \frac{27}{20})$ ,  $(0, 0)$ , and  $(3, 0)$  are critical points of the following function of three variables and classify them as local minimums, local maximums or saddle points.

$$f(x, y) = 5x^4 - 30x^3 + 45x^2 + 6yx^2 - 18yx + 5y^2.$$

18. Find the critical points of the following function of three variables and classify them as local minimums, local maximums or saddle points.

$$f(x, y, z) = \frac{10}{3}x^2 - \frac{44}{3}x + \frac{64}{3} - \frac{10}{3}yx + \frac{16}{3}y + \frac{2}{3}zx - \frac{20}{3}z + \frac{10}{3}y^2 + \frac{2}{3}zy + \frac{4}{3}z^2.$$

19. Find the critical points of the following function of three variables and classify them as local minimums, local maximums or saddle points.

$$f(x, y, z) = -\frac{7}{3}x^2 - \frac{146}{3}x + \frac{83}{3} + \frac{16}{3}yx + \frac{4}{3}y - \frac{14}{3}zx + \frac{94}{3}z - \frac{7}{3}y^2 - \frac{14}{3}zy + \frac{8}{3}z^2.$$

20. Find the critical points of the following function of three variables and classify them as local minimums, local maximums or saddle points.

$$f(x, y, z) = \frac{2}{3}x^2 + 4x + 75 - \frac{14}{3}yx - 38y - \frac{8}{3}zx - 2z + \frac{2}{3}y^2 - \frac{8}{3}zy - \frac{1}{3}z^2.$$

21. Find the critical points of the following function of three variables and classify them as local minimums, local maximums or saddle points.

$$f(x, y, z) = 4x^2 - 30x + 510 - 2yx + 60y - 2zx - 70z + 4y^2 - 2zy + 4z^2.$$

22. Show the critical points of the following function are points of the form,  $(x, y, z) = (t, 2t^2 - 10t, -t^2 + 5t)$  for  $t \in \mathbf{R}$  and classify them as local minimums, local maximums or saddle points.

$$f(x, y, z) = -\frac{1}{6}x^4 + \frac{5}{3}x^3 - \frac{25}{6}x^2 + \frac{10}{3}yx^2 - \frac{50}{3}yx + \frac{19}{3}zx^2 - \frac{95}{3}zx - \frac{5}{3}y^2 - \frac{10}{3}zy - \frac{1}{6}z^2.$$

The verification that the critical points are of the indicated form is left for you.

The Hessian is

$$\begin{pmatrix} -2x^2 + 10x - \frac{25}{3} + \frac{20}{3}y + \frac{38}{3}z & \frac{20}{3}x - \frac{50}{3} & \frac{38}{3}x - \frac{95}{3} \\ \frac{20}{3}x - \frac{50}{3} & -\frac{10}{3} & -\frac{10}{3} \\ \frac{38}{3}x - \frac{95}{3} & -\frac{10}{3} & -\frac{1}{3} \end{pmatrix}$$

at a critical point it is

$$\begin{pmatrix} -\frac{4}{3}t^2 + \frac{20}{3}t - \frac{25}{3} & \frac{20}{3}(t) - \frac{50}{3} & \frac{38}{3}(t) - \frac{95}{3} \\ \frac{20}{3}(t) - \frac{50}{3} & -\frac{10}{3} & -\frac{10}{3} \\ \frac{38}{3}(t) - \frac{95}{3} & -\frac{10}{3} & -\frac{1}{3} \end{pmatrix}.$$

The eigenvalues are

$$0, -\frac{2}{3}t^2 + \frac{10}{3}t - 6 + \frac{2}{3}\sqrt{(t^4 - 10t^3 + 493t^2 - 2340t + 2916)},$$

and

$$-\frac{2}{3}t^2 + \frac{10}{3}t - 6 - \frac{2}{3}\sqrt{(t^4 - 10t^3 + 493t^2 - 2340t + 2916)}.$$

If you graph these functions of  $t$  you find the second is always positive and the third is always negative. Therefore, all these critical points are saddle points.

23. Show the critical points of the following function are  $(0, -3, 0)$ ,  $(2, -3, 0)$ , and  $(1, -3, -\frac{1}{3})$  and classify them as local minimums, local maximums or saddle points.

$$f(x, y, z) = -\frac{3}{2}x^4 + 6x^3 - 6x^2 + zx^2 - 2zx - 2y^2 - 12y - 18 - \frac{3}{2}z^2.$$

The Hessian is

$$\begin{pmatrix} -12 + 36x + 2z - 18x^2 & 0 & -2 + 2x \\ 0 & -4 & 0 \\ -2 + 2x & 0 & -3 \end{pmatrix}$$

Now consider the critical point,  $(1, -3, -\frac{1}{3})$ . At this point the Hessian matrix equals

$$\begin{pmatrix} \frac{16}{3} & 0 & 0 \\ 0 & -4 & 0 \\ 0 & 0 & -3 \end{pmatrix},$$

The eigenvalues are  $\frac{16}{3}, -3, -4$  and so this point is a saddle point.

Next consider the critical point,  $(2, -3, 0)$ . At this point the Hessian matrix is

$$\begin{pmatrix} -12 & 0 & 2 \\ 0 & -4 & 0 \\ 2 & 0 & -3 \end{pmatrix}$$

The eigenvalues are  $-4, -\frac{15}{2} + \frac{1}{2}\sqrt{97}, -\frac{15}{2} - \frac{1}{2}\sqrt{97}$ , all negative so at this point there is a local max.

Finally consider the critical point,  $(0, -3, 0)$ . At this point the Hessian is

$$\begin{pmatrix} -12 & 0 & -2 \\ 0 & -4 & 0 \\ -2 & 0 & -3 \end{pmatrix}$$

and the eigenvalues are the same as the above, all negative. Therefore, there is a local maximum at this point.

24. Show the critical points of the following function are points of the form,  $(x, y, z) = (t, 2t^2 + 6t, -t^2 - 3t)$  for  $t \in \mathbf{R}$  and classify them as local minimums, local maximums or saddle points.

$$f(x, y, z) = -2yx^2 - 6yx - 4zx^2 - 12zx + y^2 + 2yz.$$

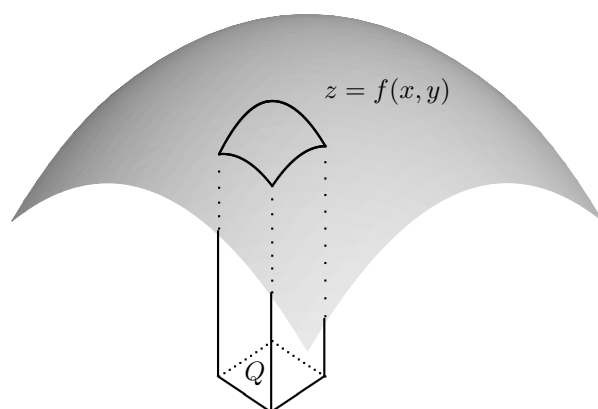
25. Show the critical points of the following function are  $(0, -1, 0)$ ,  $(4, -1, 0)$ , and  $(2, -1, -12)$  and classify them as local minimums, local maximums or saddle points.

$$f(x, y, z) = \frac{1}{2}x^4 - 4x^3 + 8x^2 - 3zx^2 + 12zx + 2y^2 + 4y + 2 + \frac{1}{2}z^2.$$

# The Riemann Integral On $\mathbb{R}^n$

## 21.1 Methods For Double Integrals

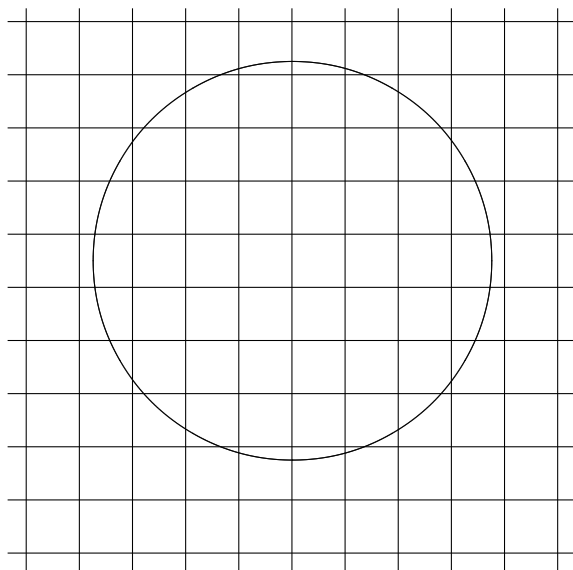
This chapter is on the Riemann integral for a function of  $n$  variables. It begins by introducing the basic concepts and applications of the integral. The proofs of the theorems involved are difficult and are left till the end. To begin with consider the problem of finding the volume under a surface of the form  $z = f(x, y)$  where  $f(x, y) \geq 0$  and  $f(x, y) = 0$  for all  $(x, y)$  outside of some bounded set. To solve this problem, consider the following picture.



In this picture, the volume of the little prism which lies above the rectangle  $Q$  and the graph of the function would lie between  $M_Q(f) v(Q)$  and  $m_Q(f) v(Q)$  where

$$M_Q(f) \equiv \sup \{f(\mathbf{x}) : \mathbf{x} \in Q\}, \quad m_Q(f) \equiv \inf \{f(\mathbf{x}) : \mathbf{x} \in Q\}, \quad (21.1)$$

and  $v(Q)$  is defined as the area of  $Q$ . Now consider the following picture.

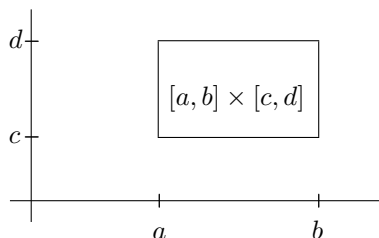


In this picture, it is assumed  $f$  equals zero outside the circle and  $f$  is a bounded nonnegative function. Then each of those little squares are the base of a prism of the sort in the previous picture and the sum of the volumes of those prisms should be the volume under the surface,  $z = f(x, y)$ . Therefore, the desired volume must lie between the two numbers,

$$\sum_Q M_Q(f) v(Q) \text{ and } \sum_Q m_Q(f) v(Q)$$

where the notation,  $\sum_Q M_Q(f) v(Q)$ , means for each  $Q$ , take  $M_Q(f)$ , multiply it by the area of  $Q$ ,  $v(Q)$ , and then add all these numbers together. Thus in  $\sum_Q M_Q(f) v(Q)$ , adds numbers which are at least as large as what is desired while in  $\sum_Q m_Q(f) v(Q)$  numbers are added which are at least as small as what is desired. Note this is a finite sum because by assumption,  $f = 0$  except for finitely many  $Q$ , namely those which intersect the circle. The sum,  $\sum_Q M_Q(f) v(Q)$  is called an upper sum,  $\sum_Q m_Q(f) v(Q)$  is a lower sum, and the desired volume is caught between these upper and lower sums.

None of this depends in any way on the function being nonnegative. It also does not depend in any essential way on the function being defined on  $\mathbb{R}^2$ , although it is impossible to draw meaningful pictures in higher dimensional cases. To define the Riemann integral, it is necessary to first give a description of something called a grid. First you must understand that something like  $[a, b] \times [c, d]$  is a rectangle in  $\mathbb{R}^2$ , having sides parallel to the axes. The situation is illustrated in the following picture.



$(x, y) \in [a, b] \times [c, d]$ , means  $x \in [a, b]$  and also  $y \in [c, d]$  and the points which do this



comprise the rectangle just as shown in the picture.

**Definition 21.1.1** For  $i = 1, 2$ , let  $\{\alpha_k^i\}_{k=-\infty}^{\infty}$  be points on  $\mathbb{R}$  which satisfy

$$\lim_{k \rightarrow \infty} \alpha_k^i = \infty, \quad \lim_{k \rightarrow -\infty} \alpha_k^i = -\infty, \quad \alpha_k^i < \alpha_{k+1}^i. \quad (21.2)$$

For such sequences, define a grid on  $\mathbb{R}^2$  denoted by  $\mathcal{G}$  or  $\mathcal{F}$  as the collection of rectangles of the form

$$Q = [\alpha_k^1, \alpha_{k+1}^1] \times [\alpha_l^2, \alpha_{l+1}^2]. \quad (21.3)$$

If  $\mathcal{G}$  is a grid, another grid,  $\mathcal{F}$  is a refinement of  $\mathcal{G}$  if every box of  $\mathcal{G}$  is the union of boxes of  $\mathcal{F}$ .

For  $\mathcal{G}$  a grid, the expression,

$$\sum_{Q \in \mathcal{G}} M_Q(f) v(Q)$$

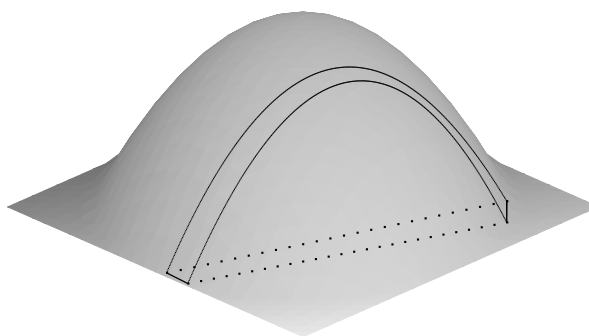
is called the upper sum associated with the grid,  $\mathcal{G}$  as described above in the discussion of the volume under a surface. Again, this means to take a rectangle from  $\mathcal{G}$  multiply  $M_Q(f)$  defined in (21.1) by its area,  $v(Q)$  and sum all these products for every  $Q \in \mathcal{G}$ . The symbol,

$$\sum_{Q \in \mathcal{G}} m_Q(f) v(Q),$$

called a lower sum, is defined similarly. With this preparation it is time to give a definition of the Riemann integral of a function of two variables.

**Definition 21.1.2** Let  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  be a bounded function which equals zero for all  $(x, y)$  outside some bounded set. Then  $\int f dV$  is defined to be the unique number which lies between all upper sums and all lower sums. In the case of  $\mathbb{R}^2$ , it is common to replace the  $V$  with  $A$  and write this symbol as  $\int f dA$  where  $A$  stands for area.

This definition begs a difficult question. For which functions does there exist a unique number between all the upper and lower sums? This interesting question is discussed in the section on the theory of the Riemann integral. First consider the question: How can the Riemann integral be computed? Consider the following picture in which  $f$  equals zero outside the rectangle  $[a, b] \times [c, d]$ .



It depicts a slice taken from the solid defined by  $\{(x, y) : 0 \leq y \leq f(x, y)\}$ . You see these when you look at a loaf of bread. If you wanted to find the volume of the loaf of bread, and you knew the volume of each slice of bread, you could find the volume of the whole loaf by adding the volumes of individual slices. It is the same here. If you could find the volume of the slice represented in this picture, you could add these up and get the volume of the solid. The slice in the picture corresponds to constant  $y$  and is assumed to be very thin, having thickness equal to  $h$ . Denote the volume of the solid under the graph of  $z = f(x, y)$  on  $[a, b] \times [c, y]$  by  $V(y)$ . Then

$$V(y+h) - V(y) \approx h \int_a^b f(x, y) dx$$

where the integral is obtained by fixing  $y$  and integrating with respect to  $x$ . It is hoped that the approximation would be increasingly good as  $h$  gets smaller. Thus, dividing by  $h$  and taking a limit, it is expected that

$$V'(y) = \int_a^b f(x, y) dx, \quad V(c) = 0.$$

Therefore, the volume of the solid under the graph of  $z = f(x, y)$  is given by

$$\int_c^d \left( \int_a^b f(x, y) dx \right) dy \quad (21.4)$$

but this was also the result of  $\int f dV$ . Therefore, it is expected that this is a way to evaluate  $\int f dV$ . Note what has been gained here. A hard problem, finding  $\int f dV$ , is reduced to a sequence of easier problems. First do

$$\int_a^b f(x, y) dx$$

getting a function of  $y$ , say  $F(y)$  and then do

$$\int_c^d \left( \int_a^b f(x, y) dx \right) dy = \int_c^d F(y) dy.$$

Of course there is nothing special about fixing  $y$  first. The same thing should be obtained from the integral,

$$\int_a^b \left( \int_c^d f(x, y) dy \right) dx \quad (21.5)$$

These expressions in (21.4) and (21.5) are called iterated integrals. They are tools for evaluating  $\int f dV$  which would be hard to find otherwise. In practice, the parenthesis is usually omitted in these expressions. Thus

$$\int_a^b \left( \int_c^d f(x, y) dy \right) dx = \int_a^b \int_c^d f(x, y) dy dx$$

and it is understood that you are to do the inside integral first and then when you have done it, obtaining a function of  $x$ , you integrate this function of  $x$ .

I have presented this for the case where  $f(x, y) \geq 0$  and the integral represents a volume, but there is no difference in the general case where  $f$  is not necessarily nonnegative.

Throughout, I have been assuming the notion of volume has some sort of independent meaning. This assumption is nonsense and is one of many reasons the above explanation does not rise to the level of a proof. It is only intended to make things plausible. A careful presentation will be given later.

Another aspect of this is the notion of integrating a function which is defined on some set, not on all  $\mathbb{R}^2$ . For example, suppose  $f$  is defined on the set,  $S \subseteq \mathbb{R}^2$ . What is meant by  $\int_S f dV$ ?

**Definition 21.1.3** Let  $f : S \rightarrow \mathbb{R}$  where  $S$  is a subset of  $\mathbb{R}^2$ . Then denote by  $f_1$  the function defined by

$$f_1(x, y) \equiv \begin{cases} f(x, y) & \text{if } (x, y) \in S \\ 0 & \text{if } (x, y) \notin S \end{cases}.$$

Then

$$\int_S f dV \equiv \int f_1 dV.$$

**Example 21.1.4** Let  $f(x, y) = x^2y + yx$  for  $(x, y) \in [0, 1] \times [0, 2] \equiv R$ . Find  $\int_R f dV$ .

This is done using iterated integrals like those defined above. Thus

$$\int_R f dV = \int_0^1 \int_0^2 (x^2y + yx) dy dx.$$

The inside integral yields

$$\int_0^2 (x^2y + yx) dy = 2x^2 + 2x$$

and now the process is completed by doing  $\int_0^1$  to what was just obtained. Thus

$$\int_0^1 \int_0^2 (x^2y + yx) dy dx = \int_0^1 (2x^2 + 2x) dx = \frac{5}{3}.$$

If the integration is done in the opposite order, the same answer should be obtained.

$$\int_0^2 \int_0^1 (x^2y + yx) dx dy$$

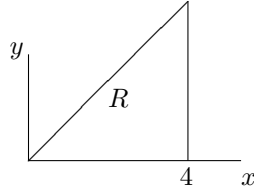
$$\int_0^1 (x^2y + yx) dx = \frac{5}{6}y$$

Now

$$\int_0^2 \int_0^1 (x^2y + yx) dx dy = \int_0^2 \left(\frac{5}{6}y\right) dy = \frac{5}{3}.$$

If a different answer had been obtained it would have been a sign that a mistake had been made.

**Example 21.1.5** Let  $f(x, y) = x^2y + yx$  for  $(x, y) \in R$  where  $R$  is the triangular region defined to be in the first quadrant, below the line  $y = x$  and to the left of the line  $x = 4$ . Find  $\int_R f dV$ .



Now from the above discussion,

$$\int_R f dV = \int_0^4 \int_0^x (x^2 y + yx) dy dx$$

The reason for this is that  $x$  goes from 0 to 4 and for each fixed  $x$  between 0 and 4,  $y$  goes from 0 to the slanted line,  $y = x$ . Thus  $y$  goes from 0 to  $x$ . This explains the inside integral. Now  $\int_0^x (x^2 y + yx) dy = \frac{1}{2}x^4 + \frac{1}{2}x^3$  and so

$$\int_R f dV = \int_0^4 \left( \frac{1}{2}x^4 + \frac{1}{2}x^3 \right) dx = \frac{672}{5}.$$

What of integration in a different order? Lets put the integral with respect to  $y$  on the outside and the integral with respect to  $x$  on the inside. Then

$$\int_R f dV = \int_0^4 \int_y^4 (x^2 y + yx) dx dy$$

For each  $y$  between 0 and 4, the variable  $x$ , goes from  $y$  to 4.

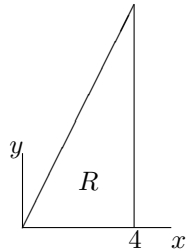
$$\int_y^4 (x^2 y + yx) dx = \frac{88}{3}y - \frac{1}{3}y^4 - \frac{1}{2}y^3$$

Now

$$\int_R f dV = \int_0^4 \left( \frac{88}{3}y - \frac{1}{3}y^4 - \frac{1}{2}y^3 \right) dy = \frac{672}{5}.$$

Here is a similar example.

**Example 21.1.6** Let  $f(x, y) = x^2 y$  for  $(x, y) \in R$  where  $R$  is the triangular region defined to be in the first quadrant, below the line  $y = 2x$  and to the left of the line  $x = 4$ . Find  $\int_R f dV$ .



Put the integral with respect to  $x$  on the outside first. Then

$$\int_R f dV = \int_0^4 \int_0^{2x} (x^2 y) dy dx$$

because for each  $x \in [0, 4]$ ,  $y$  goes from 0 to  $2x$ . Then

$$\int_0^{2x} (x^2 y) dy = 2x^4$$

and so

$$\int_R f dV = \int_0^4 (2x^4) dx = \frac{2048}{5}$$

Now do the integral in the other order. Here the integral with respect to  $y$  will be on the outside. What are the limits of this integral? Look at the triangle and note that  $x$  goes from 0 to 4 and so  $2x = y$  goes from 0 to 8. Now for fixed  $y$  between 0 and 8, where does  $x$  go? It goes from the  $x$  coordinate on the line  $y = 2x$  which corresponds to this  $y$  to 4. What is the  $x$  coordinate on this line which goes with  $y$ ? It is  $x = y/2$ . Therefore, the iterated integral is

$$\int_0^8 \int_{y/2}^4 (x^2 y) dx dy.$$

Now

$$\int_{y/2}^4 (x^2 y) dx = \frac{64}{3} y - \frac{1}{24} y^4$$

and so

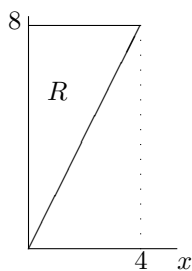
$$\int_R f dV = \int_0^8 \left( \frac{64}{3} y - \frac{1}{24} y^4 \right) dy = \frac{2048}{5}$$

the same answer.

A few observations are in order here. In finding  $\int_S f dV$  there is no problem in setting things up if  $S$  is a rectangle. However, if  $S$  is not a rectangle, the procedure **always** is agonizing. A good rule of thumb is that if what you do is easy it will be wrong. There are no shortcuts! There are no quick fixes which require no thought! Pain and suffering is inevitable and you must not expect it to be otherwise. Always draw a picture and then begin agonizing over the correct limits. Even when you are careful you will make lots of mistakes until you get used to the process.

Sometimes an integral can be evaluated in one order but not in another.

**Example 21.1.7** For  $R$  as shown below, find  $\int_R \sin(y^2) dV$ .



Setting this up to have the integral with respect to  $y$  on the inside yields

$$\int_0^4 \int_{2x}^8 \sin(y^2) dy dx.$$

Unfortunately, there is no antiderivative in terms of elementary functions for  $\sin(y^2)$  so there is an immediate problem in evaluating the inside integral. It doesn't work out so the

next step is to do the integration in another order and see if some progress can be made. This yields

$$\int_0^8 \int_0^{y/2} \sin(y^2) \, dx \, dy = \int_0^8 \frac{y}{2} \sin(y^2) \, dy$$

and  $\int_0^8 \frac{y}{2} \sin(y^2) \, dy = -\frac{1}{4} \cos 64 + \frac{1}{4}$  which you can verify by making the substitution,  $u = y^2$ . Thus

$$\int_R \sin(y^2) \, dy = -\frac{1}{4} \cos 64 + \frac{1}{4}.$$

This illustrates an important idea. The integral  $\int_R \sin(y^2) \, dV$  is defined as a number. It is the unique number between all the upper sums and all the lower sums. Finding it is another matter. In this case it was possible to find it using one order of integration but not the other. The iterated integral in this other order also is defined as a number but it can't be found directly without interchanging the order of integration. Of course sometimes nothing you try will work out.

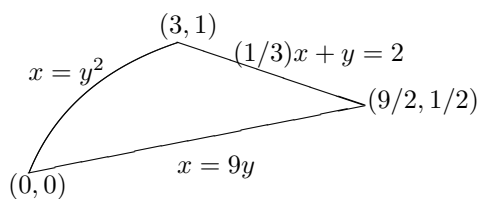
Consider a two dimensional material. Of course there is no such thing but a flat plate might be modeled as one. The density  $\rho$  is a function of position and is defined as follows. Consider a small chunk of area,  $dV$  located at the point whose Cartesian coordinates are  $(x, y)$ . Then the mass of this small chunk of material is given by  $\rho(x, y) \, dV$ . Thus if the material occupies a region in two dimensional space,  $U$ , the total mass of this material would be

$$\int_U \rho \, dV$$

In other words you integrate the density to get the mass. Now by letting  $\rho$  depend on position, you can include the case where the material is not homogeneous. Here is an example.

**Example 21.1.8** Let  $\rho(x, y)$  denote the density of the plane region determined by the curves  $\frac{1}{3}x + y = 2$ ,  $x = 3y^2$ , and  $x = 9y$ . Find the total mass if  $\rho(x, y) = y$ .

You need to first draw a picture of the region,  $R$ . A rough sketch follows.



This region is in two pieces, one having the graph of  $x = 9y$  on the bottom and the graph of  $x = 3y^2$  on the top and another piece having the graph of  $x = 9y$  on the bottom and the graph of  $\frac{1}{3}x + y = 2$  on the top. Therefore, in setting up the integrals, with the integral with respect to  $x$  on the outside, the double integral equals the following sum of iterated integrals.

$$\overbrace{\int_0^3 \int_{x/9}^{\sqrt{x/3}} y \, dy \, dx}^{\text{has } x=3y^2 \text{ on top}} + \overbrace{\int_3^{9/2} \int_{x/9}^{2-\frac{1}{3}x} y \, dy \, dx}^{\text{has } \frac{1}{3}x+y=2 \text{ on top}}$$

You notice it is not necessary to have a perfect picture, just one which is good enough to figure out what the limits should be. The dividing line between the two cases is  $x = 3$  and this was shown in the picture. Now it is only a matter of evaluating the iterated integrals which in this case is routine and gives 1.

## 21.2 Exercises

1. Let  $\rho(x, y)$  denote the density of the plane region determined by the curves  $\frac{1}{4}x + y = 6$ ,  $x = 4y^2$ , and  $x = 16y$ . Find the total mass if  $\rho(x, y) = y$ . Your answer should be  $\frac{1168}{75}$ .
2. Let  $\rho(x, y)$  denote the density of the plane region determined by the curves  $\frac{1}{5}x + y = 6$ ,  $x = 5y^2$ , and  $x = 25y$ . Find the total mass if  $\rho(x, y) = y + 2x$ . Your answer should be  $\frac{1735}{3}$ .
3. Let  $\rho(x, y)$  denote the density of the plane region determined by the curves  $y = 3x$ ,  $y = x$ ,  $3x + 3y = 9$ . Find the total mass if  $\rho(x, y) = y + 1$ . Your answer should be  $\frac{81}{32}$ .
4. Let  $\rho(x, y)$  denote the density of the plane region determined by the curves  $y = 3x$ ,  $y = x$ ,  $4x + 2y = 8$ . Find the total mass if  $\rho(x, y) = y + 1$ .
5. Let  $\rho(x, y)$  denote the density of the plane region determined by the curves  $y = 3x$ ,  $y = x$ ,  $2x + 2y = 4$ . Find the total mass if  $\rho(x, y) = x + 2y$ .
6. Let  $\rho(x, y)$  denote the density of the plane region determined by the curves  $y = 3x$ ,  $y = x$ ,  $5x + 2y = 10$ . Find the total mass if  $\rho(x, y) = y + 1$ .
7. Find  $\int_0^4 \int_{y/2}^2 \frac{1}{x} e^{2\frac{y}{x}} dx dy$ . Your answer should be  $e^4 - 1$ . You might need to interchange the order of integration.
8. Find  $\int_0^8 \int_{y/2}^4 \frac{1}{x} e^{3\frac{y}{x}} dx dy$ .
9. Find  $\int_0^8 \int_{y/2}^4 \frac{1}{x} e^{3\frac{y}{x}} dx dy$ .
10. Find  $\int_0^4 \int_{y/2}^2 \frac{1}{x} e^{3\frac{y}{x}} dx dy$ .
11. Evaluate the iterated integral and then write the iterated integral with the order of integration reversed.  $\int_0^4 \int_0^{3y} xy^3 dx dy$ . Your answer for the iterated integral should be  $\int_0^{12} \int_{\frac{1}{3}x}^4 xy^3 dy dx$ .
12. Evaluate the iterated integral and then write the iterated integral with the order of integration reversed.  $\int_0^3 \int_0^{3y} xy^3 dx dy$ .
13. Evaluate the iterated integral and then write the iterated integral with the order of integration reversed.  $\int_0^2 \int_0^{2y} xy^2 dx dy$ .
14. Evaluate the iterated integral and then write the iterated integral with the order of integration reversed.  $\int_0^3 \int_0^y xy^3 dx dy$ .
15. Evaluate the iterated integral and then write the iterated integral with the order of integration reversed.  $\int_0^1 \int_0^y xy^2 dx dy$ .
16. Evaluate the iterated integral and then write the iterated integral with the order of integration reversed.  $\int_0^5 \int_0^{3y} xy^2 dx dy$ .

17. Find  $\int_0^{\frac{1}{3}\pi} \int_x^{\frac{1}{3}\pi} \frac{\sin y}{y} dy dx$ . Your answer should be  $\frac{1}{2}$ .
18. Find  $\int_0^{\frac{1}{2}\pi} \int_x^{\frac{1}{2}\pi} \frac{\sin y}{y} dy dx$ .
19. Find  $\int_0^\pi \int_x^\pi \frac{\sin y}{y} dy dx$ .
20. Evaluate the iterated integral and then write the iterated integral with the order of integration reversed.  $\int_{-3}^3 \int_{-x}^x x^2 dy dx$
- Your answer for the iterated integral should be  $\int_3^0 \int_{-3}^{-y} x^2 dx dy + \int_0^{-3} \int_{-3}^y x^2 dx dy + \int_0^3 \int_y^3 x^2 dx dy + \int_{-3}^0 \int_{-y}^3 x^2 dx dy$ . This is a very interesting example which shows that iterated integrals have a life of their own, not just as a method for evaluating double integrals.
21. Evaluate the iterated integral and then write the iterated integral with the order of integration reversed.  $\int_{-2}^2 \int_{-x}^x x^2 dy dx$ .

## 21.3 Methods For Triple Integrals

The integral of a function of three variables is defined similar to the integral of a function of two variables.

**Definition 21.3.1** For  $i = 1, 2, 3$  let  $\{\alpha_k^i\}_{k=-\infty}^\infty$  be points on  $\mathbb{R}$  which satisfy

$$\lim_{k \rightarrow \infty} \alpha_k^i = \infty, \quad \lim_{k \rightarrow -\infty} \alpha_k^i = -\infty, \quad \alpha_k^i < \alpha_{k+1}^i. \quad (21.6)$$

For such sequences, define a grid on  $\mathbb{R}^3$  denoted by  $\mathcal{G}$  or  $\mathcal{F}$  as the collection of boxes of the form

$$Q = [\alpha_k^1, \alpha_{k+1}^1] \times [\alpha_l^2, \alpha_{l+1}^2] \times [\alpha_p^3, \alpha_{p+1}^3]. \quad (21.7)$$

If  $\mathcal{G}$  is a grid,  $\mathcal{F}$  is called a refinement of  $\mathcal{G}$  if every box of  $\mathcal{G}$  is the union of boxes of  $\mathcal{F}$ .

For  $\mathcal{G}$  a grid,

$$\sum_{Q \in \mathcal{G}} M_Q(f) v(Q)$$

is the upper sum associated with the grid,  $\mathcal{G}$  where

$$M_Q(f) \equiv \sup \{f(\mathbf{x}) : \mathbf{x} \in Q\}$$

and if  $Q = [a, b] \times [c, d] \times [e, f]$ , then  $v(Q)$  is the volume of  $Q$  given by  $(b-a)(d-c)(f-e)$ . Letting

$$m_Q(f) \equiv \inf \{f(\mathbf{x}) : \mathbf{x} \in Q\}$$

the lower sum associated with this partition is

$$\sum_{Q \in \mathcal{G}} m_Q(f) v(Q),$$

With this preparation it is time to give a definition of the Riemann integral of a function of three variables. This definition is just like the one for a function of two variables.

**Definition 21.3.2** Let  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  be a bounded function which equals zero outside of some bounded subset of  $\mathbb{R}^3$ .  $\int f dV$  is defined as the unique number between all the upper sums and lower sums.



As in the case of a function of two variables there are all sorts of mathematical questions which are dealt with later.

The way to think of integrals is as follows. Located at a point  $\mathbf{x}$ , there is an “infinitesimal” chunk of volume,  $dV$ . The integral involves taking this little chunk of volume,  $dV$ , multiplying it by  $f(\mathbf{x})$  and then adding up all such products. Upper sums are too large and lower sums are too small but the unique number between all the lower and upper sums is just right and corresponds to the notion of adding up all the  $f(\mathbf{x}) dV$ . Even the notation is suggestive of this concept of sum. It is a long thin  $S$  denoting sum. This is the fundamental concept for the integral in any number of dimensions and all the definitions and technicalities are designed to give precision and mathematical respectability to this notion.

To consider how to evaluate triple integrals, imagine a sum of the form  $\sum_{ijk} a_{ijk}$  where there are only finitely many choices for  $i, j$ , and  $k$  and the symbol means you simply add up all the  $a_{ijk}$ . By the commutative law of addition, these may be added systematically in the form,  $\sum_k \sum_j \sum_i a_{ijk}$ . A similar process is used to evaluate triple integrals and since integrals are like sums, you might expect it to be valid. Specifically,

$$\int f dV = \int \int \int f(x, y, z) dx dy dz.$$

In words, sum with respect to  $x$  and then sum what you get with respect to  $y$  and finally, with respect to  $z$ . Of course this should hold in any other order such as

$$\int f dV = \int \int \int f(x, y, z) dz dy dx.$$

Later this will be proved under appropriate conditions<sup>1</sup>.

Having discussed double and triple integrals, the definition of the integral of a function of  $n$  variables is accomplished in the same way. It is given in Definition 21.11.3 on Page 559. In this book most integrals will involve no more than three variables. However, this does not mean an integral of a function of more than three variables is unimportant. Therefore, I will begin to refer to the general case when theorems are stated.

**Definition 21.3.3** For  $E \subseteq \mathbb{R}^n$ ,

$$\mathcal{X}_E(\mathbf{x}) \equiv \begin{cases} 1 & \text{if } \mathbf{x} \in E \\ 0 & \text{if } \mathbf{x} \notin E \end{cases}.$$

Define  $\int_E f dV \equiv \int \mathcal{X}_E f dV$  when  $f \mathcal{X}_E \in \mathcal{R}(\mathbb{R}^n)$ .

As an example of the use of triple integrals, consider a solid occupying a set of points,  $U \subseteq \mathbb{R}^3$  having density  $\rho$ . Thus  $\rho$  is a function of position and the total mass of the solid equals

$$\int_U \rho dV.$$

This is just like the two dimensional case. The mass of an infinitesimal chunk of the solid located at  $\mathbf{x}$  would be  $\rho(\mathbf{x}) dV$  and so the total mass is just the sum of all these,  $\int_U \rho(\mathbf{x}) dV$ .

**Example 21.3.4** Find  $\int_2^3 \int_3^x \int_{3y}^x (x - y) dz dy dx$ .

---

<sup>1</sup>All of these fundamental questions about integrals can be considered more easily in the context of the Lebesgue integral. However, this integral is more abstract than the Riemann integral.

The inside integral yields  $\int_{3y}^x (x - y) dz = x^2 - 4xy + 3y^2$ . Next this must be integrated with respect to  $y$  to give  $\int_3^x (x^2 - 4xy + 3y^2) dy = -3x^2 + 18x - 27$ . Finally the third integral gives

$$\int_2^3 \int_3^x \int_{3y}^x (x - y) dz dy dx = \int_2^3 (-3x^2 + 18x - 27) dx = -1.$$

**Example 21.3.5** Find  $\int_0^\pi \int_0^{3y} \int_0^{y+z} \cos(x + y) dx dz dy$ .

The inside integral is  $\int_0^{y+z} \cos(x + y) dx = 2 \cos z \sin y \cos y + 2 \sin z \cos^2 y - \sin z - \sin y$ . Now this has to be integrated.

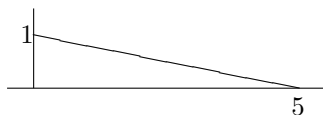
$$\begin{aligned} \int_0^{3y} \int_0^{y+z} \cos(x + y) dx dz &= \int_0^{3y} (2 \cos z \sin y \cos y + 2 \sin z \cos^2 y - \sin z - \sin y) dz \\ &= -1 - 16 \cos^5 y + 20 \cos^3 y - 5 \cos y - 3(\sin y)y + 2 \cos^2 y. \end{aligned}$$

Finally, this last expression must be integrated from 0 to  $\pi$ . Thus

$$\begin{aligned} &\int_0^\pi \int_0^{3y} \int_0^{y+z} \cos(x + y) dx dz dy \\ &= \int_0^\pi (-1 - 16 \cos^5 y + 20 \cos^3 y - 5 \cos y - 3(\sin y)y + 2 \cos^2 y) dy \\ &= -3\pi \end{aligned}$$

**Example 21.3.6** Find the volume of  $R$  where  $R$  is the bounded region formed by the plane  $\frac{1}{5}x + y + \frac{1}{5}z = 1$  and the planes  $x = 0, y = 0, z = 0$ .

When  $z = 0$ , the plane becomes  $\frac{1}{5}x + y = 1$ . Thus the intersection of this plane with the  $xy$  plane is this line shown in the following picture.



Therefore, the bounded region is between the triangle formed in the above picture by the  $x$  axis, the  $y$  axis and the above line and the surface given by  $\frac{1}{5}x + y + \frac{1}{5}z = 1$  or  $z = 5(1 - (\frac{1}{5}x + y)) = 5 - x - 5y$ . Therefore, an iterated integral which yields the volume is

$$\int_0^5 \int_0^{1-\frac{1}{5}x} \int_0^{5-x-5y} dz dy dx = \frac{25}{6}.$$

**Example 21.3.7** Find the mass of the bounded region,  $R$  formed by the plane  $\frac{1}{3}x + \frac{1}{3}y + \frac{1}{5}z = 1$  and the planes  $x = 0, y = 0, z = 0$  if the density is  $\rho(x, y, z) = z$ .

This is done just like the previous example except in this case there is a function to integrate. Thus the answer is

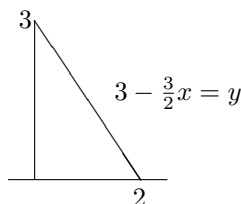
$$\int_0^3 \int_0^{3-x} \int_0^{5-\frac{5}{3}x-\frac{5}{3}y} z dz dy dx = \frac{75}{8}.$$

**Example 21.3.8** Here is an iterated integral:  $\int_0^2 \int_0^{3-\frac{3}{2}x} \int_0^{x^2} dz dy dx$ . Write as an iterated integral in the order  $dz dx dy$ .

The inside integral is just a function of  $x$  and  $y$ . (In fact, only a function of  $x$ .) The order of the last two integrals must be interchanged. Thus the iterated integral which needs to be done in a different order is

$$\int_0^2 \int_0^{3-\frac{3}{2}x} f(x, y) dy dx.$$

As usual, it is important to draw a picture and then go from there.



Thus this double integral equals

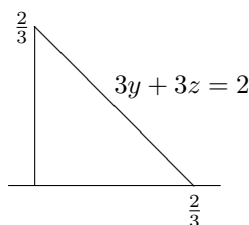
$$\int_0^3 \int_0^{\frac{2}{3}(3-y)} f(x, y) dx dy.$$

Now substituting in for  $f(x, y)$ ,

$$\int_0^3 \int_0^{\frac{2}{3}(3-y)} \int_0^{x^2} dz dx dy.$$

**Example 21.3.9** Find the volume of the bounded region determined by  $3y + 3z = 2$ ,  $x = 16 - y^2$ ,  $y = 0$ ,  $x = 0$ .

In the  $yz$  plane, the following picture corresponds to  $x = 0$ .



Therefore, the outside integrals taken with respect to  $z$  and  $y$  are of the form  $\int_0^{\frac{2}{3}} \int_0^{\frac{2}{3}-y} dz dy$  and now for any choice of  $(y, z)$  in the above triangular region,  $x$  goes from 0 to  $16 - y^2$ . Therefore, the iterated integral is

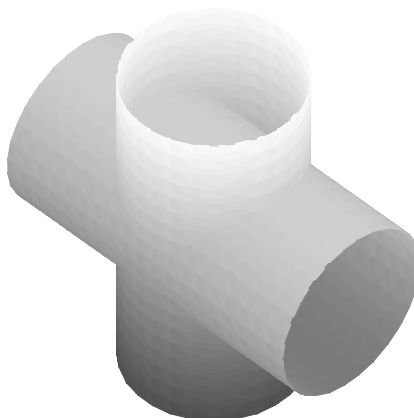
$$\int_0^{\frac{2}{3}} \int_0^{\frac{2}{3}-y} \int_0^{16-y^2} dx dz dy = \frac{860}{243}$$

**Example 21.3.10** Find the volume of the region determined by the intersection of the two cylinders,  $x^2 + y^2 \leq 9$  and  $y^2 + z^2 \leq 9$ .

The first listed cylinder intersects the  $xy$  plane in the disk,  $x^2 + y^2 \leq 9$ . What is the volume of the three dimensional region which is between this disk and the two surfaces,  $z = \sqrt{9 - y^2}$  and  $z = -\sqrt{9 - y^2}$ ? An iterated integral for the volume is

$$\int_{-3}^3 \int_{-\sqrt{9-y^2}}^{\sqrt{9-y^2}} \int_{-\sqrt{9-y^2}}^{\sqrt{9-y^2}} dz dx dy = 144.$$

Note I drew no picture of the three dimensional region. If you are interested, here it is.



One of the cylinders is parallel to the  $z$  axis,  $x^2 + y^2 \leq 9$  and the other is parallel to the  $x$  axis,  $y^2 + z^2 \leq 9$ . I did not need to be able to draw such a nice picture in order to work this problem. This is the key to doing these. Draw pictures in two dimensions and reason from the two dimensional pictures rather than attempt to wax artistic and consider all three dimensions at once. These problems are hard enough without making them even harder by attempting to be an artist.

**Example 21.3.11** Find the total mass of the bounded solid determined by  $z = 9 - x^2 - y^2$  and  $x, y, z \geq 0$  if the mass is given by  $\rho(x, y, z) = z$

When  $z = 0$  the surface,  $z = 9 - x^2 - y^2$  intersects the  $xy$  plane in a circle of radius 3 centered at  $(0, 0)$ . Since  $x, y \geq 0$ , it is only a quarter of a circle of interest, the part where both these variables are nonnegative. For each  $(x, y)$  inside this quarter circle,  $z$  goes from 0 to  $9 - x^2 - y^2$ . Therefore, the iterated integral is of the form,

$$\int_0^3 \int_0^{\sqrt{9-x^2}} \int_0^{9-x^2-y^2} z \, dz \, dy \, dx = \frac{243}{8}\pi$$

**Example 21.3.12** Find the volume of the bounded region determined by  $x \geq 0, y \geq 0, z \geq 0$ , and  $\frac{1}{7}x + y + \frac{1}{4}z = 1$ , and  $x + \frac{1}{7}y + \frac{1}{4}z = 1$ .

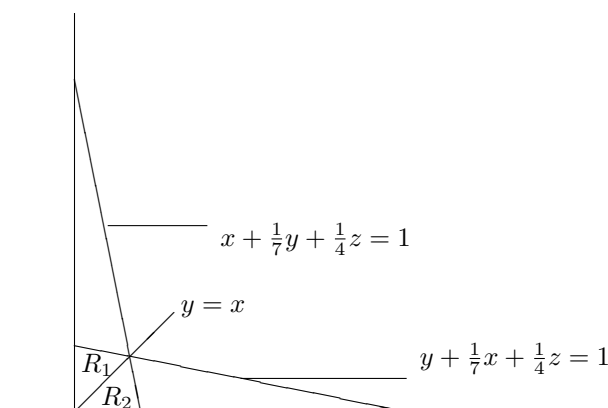
When  $z = 0$ , the plane  $\frac{1}{7}x + y + \frac{1}{4}z = 1$  intersects the  $xy$  plane in the line whose equation is

$$\frac{1}{7}x + y = 1$$

while the plane,  $x + \frac{1}{7}y + \frac{1}{4}z = 1$  intersects the  $xy$  plane in the line whose equation is

$$x + \frac{1}{7}y = 1.$$

Furthermore, the two planes intersect when  $x = y$  as can be seen from the equations,  $x + \frac{1}{7}y = 1 - \frac{z}{4}$  and  $\frac{1}{7}x + y = 1 - \frac{z}{4}$  which imply  $x = y$ . Thus the two dimensional picture to look at is depicted in the following picture.



You see in this picture, the base of the region in the  $xy$  plane is the union of the two triangles,  $R_1$  and  $R_2$ . For  $(x, y) \in R_1$ ,  $z$  goes from 0 to what it needs to be to be on the plane,  $\frac{1}{7}x + y + \frac{1}{4}z = 1$ . Thus  $z$  goes from 0 to  $4(1 - \frac{1}{7}x - y)$ . Similarly, on  $R_2$ ,  $z$  goes from 0 to  $4(1 - \frac{1}{7}y - x)$ . Therefore, the integral needed is

$$\int_{R_1} \int_0^{4(1-\frac{1}{7}x-y)} dz dV + \int_{R_2} \int_0^{4(1-\frac{1}{7}y-x)} dz dV$$

and now it only remains to consider  $\int_{R_1} dV$  and  $\int_{R_2} dV$ . The point of intersection of these lines shown in the above picture is  $(\frac{7}{8}, \frac{7}{8})$  and so an iterated integral is

$$\int_0^{7/8} \int_x^{1-\frac{x}{7}} \int_0^{4(1-\frac{1}{7}x-y)} dz dy dx + \int_0^{7/8} \int_y^{1-\frac{y}{7}} \int_0^{4(1-\frac{1}{7}x-y)} dz dx dy = \frac{7}{6}.$$

## 21.4 Exercises With Answers

The evaluation of integrals by setting up appropriate iterated integrals and then evaluating these requires a lot of practice. Therefore, I have included exercises with answers. Each of these exercises corresponds to one which does not have answers in the next section.

1. Evaluate the integral  $\int_4^7 \int_5^{3x} \int_{5y}^x dz dy dx$

Answer:

$$-\frac{3417}{2}$$

2. Find  $\int_0^4 \int_0^{2-5x} \int_0^{4-2x-y} (2x) dz dy dx$

Answer:

$$-\frac{2464}{3}$$

3. Find  $\int_0^2 \int_0^{2-5x} \int_0^{1-4x-3y} (2x) dz dy dx$

Answer:

$$-\frac{196}{3}$$

4. Evaluate the integral  $\int_5^8 \int_4^{3x} \int_{4y}^x (x - y) dz dy dx$

Answer:

$$\frac{114607}{8}$$

5. Evaluate the integral  $\int_0^\pi \int_0^{4y} \int_0^{y+z} \cos(x+y) \, dx \, dz \, dy$

Answer:

$$-4\pi$$

6. Evaluate the integral  $\int_0^\pi \int_0^{2y} \int_0^{y+z} \sin(x+y) \, dx \, dz \, dy$

Answer:

$$-\frac{19}{4}$$

7. Fill in the missing limits.  $\int_0^1 \int_0^z \int_0^z f(x, y, z) \, dx \, dy \, dz = \int_?^? \int_?^? \int_?^? f(x, y, z) \, dx \, dz \, dy,$

$$\int_0^1 \int_0^z \int_0^{2z} f(x, y, z) \, dx \, dy \, dz = \int_?^? \int_?^? \int_?^? f(x, y, z) \, dy \, dz \, dx,$$

$$\int_0^1 \int_0^z \int_0^z f(x, y, z) \, dx \, dy \, dz = \int_?^? \int_?^? \int_?^? f(x, y, z) \, dz \, dy \, dx,$$

$$\int_0^1 \int_{z/2}^{\sqrt{z}} \int_0^{y+z} f(x, y, z) \, dx \, dy \, dz = \int_?^? \int_?^? \int_?^? f(x, y, z) \, dx \, dz \, dy,$$

$$\int_5^7 \int_2^5 \int_0^3 f(x, y, z) \, dx \, dy \, dz = \int_?^? \int_?^? \int_?^? f(x, y, z) \, dz \, dy \, dx.$$

Answer:

$$\int_0^1 \int_0^z \int_0^z f(x, y, z) \, dx \, dy \, dz = \int_0^1 \int_y^1 \int_0^z f(x, y, z) \, dx \, dz \, dy,$$

$$\int_0^1 \int_0^z \int_0^{2z} f(x, y, z) \, dx \, dy \, dz = \int_0^2 \int_{x/2}^1 \int_0^z f(x, y, z) \, dy \, dz \, dx,$$

$$\int_0^1 \int_0^z \int_0^z f(x, y, z) \, dx \, dy \, dz = \int_0^1 \left[ \int_0^x \int_x^1 f(x, y, z) \, dz \, dy + \int_x^1 \int_y^1 f(x, y, z) \, dz \, dy \right] dx,$$

$$\int_0^1 \int_{z/2}^{\sqrt{z}} \int_0^{y+z} f(x, y, z) \, dx \, dy \, dz =$$

$$\int_0^{1/2} \int_{y^2}^{2y} \int_0^{y+z} f(x, y, z) \, dx \, dz \, dy + \int_{1/2}^1 \int_{y^2}^1 \int_0^{y+z} f(x, y, z) \, dx \, dz \, dy$$

$$\int_5^7 \int_2^5 \int_0^3 f(x, y, z) \, dx \, dy \, dz = \int_0^3 \int_2^5 \int_5^7 f(x, y, z) \, dz \, dy \, dx$$

8. Find the volume of  $R$  where  $R$  is the bounded region formed by the plane  $\frac{1}{5}x + y + \frac{1}{4}z = 1$  and the planes  $x = 0, y = 0, z = 0$ .

Answer:

$$\int_0^5 \int_0^{1-\frac{1}{5}x} \int_0^{4-\frac{4}{5}x-4y} dz \, dy \, dx = \frac{10}{3}$$

9. Find the volume of  $R$  where  $R$  is the bounded region formed by the plane  $\frac{1}{5}x + \frac{1}{2}y + \frac{1}{4}z = 1$  and the planes  $x = 0, y = 0, z = 0$ .

Answer:

$$\int_0^5 \int_0^{2-\frac{2}{5}x} \int_0^{4-\frac{4}{5}x-2y} dz \, dy \, dx = \frac{20}{3}$$

10. Find the mass of the bounded region,  $R$  formed by the plane  $\frac{1}{4}x + \frac{1}{2}y + \frac{1}{3}z = 1$  and the planes  $x = 0, y = 0, z = 0$  if the density is  $\rho(x, y, z) = y$

Answer:

$$\int_0^4 \int_0^{2-\frac{1}{2}x} \int_0^{3-\frac{3}{4}x-\frac{3}{2}y} (y) \, dz \, dy \, dx = 2$$

11. Find the mass of the bounded region,  $R$  formed by the plane  $\frac{1}{2}x + \frac{1}{2}y + \frac{1}{4}z = 1$  and the planes  $x = 0, y = 0, z = 0$  if the density is  $\rho(x, y, z) = z^2$

Answer:

$$\int_0^2 \int_0^{2-x} \int_0^{4-2x-2y} (z^2) \, dz \, dy \, dx = \frac{64}{15}$$

12. Here is an iterated integral:  $\int_0^3 \int_0^{3-x} \int_0^{x^2} dz dy dx$ . Write as an iterated integral in the following orders:  $dz dx dy$ ,  $dx dz dy$ ,  $dx dy dz$ ,  $dy dx dz$ ,  $dy dz dx$ .

Answer:  $\int_0^3 \int_0^{x^2} \int_0^{3-x} dy dz dx$ ,  $\int_0^9 \int_{\sqrt{z}}^3 \int_0^{3-x} dy dx dz$ ,  $\int_0^9 \int_0^{3-\sqrt{z}} \int_{\sqrt{z}}^{3-y} dx dy dz$ ,  $\int_0^3 \int_0^{3-y} \int_0^{x^2} dz dx dy$ ,  $\int_0^3 \int_0^{(3-y)^2} \int_{\sqrt{z}}^{3-y} dx dz dy$

13. Find the volume of the bounded region determined by  $5y + 2z = 4$ ,  $x = 4 - y^2$ ,  $y = 0$ ,  $x = 0$ .

Answer:  
 $\int_0^{\frac{4}{5}} \int_0^{2-\frac{5}{2}y} \int_0^{4-y^2} dx dz dy = \frac{1168}{375}$

14. Find the volume of the bounded region determined by  $4y + 3z = 3$ ,  $x = 4 - y^2$ ,  $y = 0$ ,  $x = 0$ .

Answer:  
 $\int_0^{\frac{3}{4}} \int_0^{1-\frac{4}{3}y} \int_0^{4-y^2} dx dz dy = \frac{375}{256}$

15. Find the volume of the bounded region determined by  $3y + z = 3$ ,  $x = 4 - y^2$ ,  $y = 0$ ,  $x = 0$ .

Answer:  
 $\int_0^1 \int_0^{3-3y} \int_0^{4-y^2} dx dz dy = \frac{23}{4}$

16. Find the volume of the region bounded by  $x^2 + y^2 = 16$ ,  $z = 3x$ ,  $z = 0$ , and  $x \geq 0$ .

Answer:  
 $\int_0^4 \int_{-\sqrt{16-x^2}}^{\sqrt{16-x^2}} \int_0^{3x} dz dy dx = 128$

17. Find the volume of the region bounded by  $x^2 + y^2 = 25$ ,  $z = 2x$ ,  $z = 0$ , and  $x \geq 0$ .

Answer:  
 $\int_0^5 \int_{-\sqrt{25-x^2}}^{\sqrt{25-x^2}} \int_0^{2x} dz dy dx = \frac{500}{3}$

18. Find the volume of the region determined by the intersection of the two cylinders,  $x^2 + y^2 \leq 9$  and  $y^2 + z^2 \leq 9$ .

Answer:  
 $8 \int_0^3 \int_0^{\sqrt{9-y^2}} \int_0^{\sqrt{9-y^2}} dz dx dy = 144$

19. Find the total mass of the bounded solid determined by  $z = a^2 - x^2 - y^2$  and  $x, y, z \geq 0$  if the mass is given by  $\rho(x, y, z) = z$

Answer:  
 $\int_0^4 \int_0^{\sqrt{16-x^2}} \int_0^{16-x^2-y^2} (z) dz dy dx = \frac{512}{3}\pi$

20. Find the total mass of the bounded solid determined by  $z = a^2 - x^2 - y^2$  and  $x, y, z \geq 0$  if the mass is given by  $\rho(x, y, z) = x + 1$

Answer:  
 $\int_0^5 \int_0^{\sqrt{25-x^2}} \int_0^{25-x^2-y^2} (x+1) dz dy dx = \frac{625}{8}\pi + \frac{1250}{3}$

21. Find the volume of the region bounded by  $x^2 + y^2 = 9$ ,  $z = 0$ ,  $z = 5 - y$

Answer:  
 $\int_{-3}^3 \int_{-\sqrt{9-x^2}}^{\sqrt{9-x^2}} \int_0^{5-y} dz dy dx = 45\pi$

22. Find the volume of the bounded region determined by  $x \geq 0, y \geq 0, z \geq 0$ , and  $\frac{1}{2}x + y + \frac{1}{2}z = 1$ , and  $x + \frac{1}{2}y + \frac{1}{2}z = 1$ .

Answer:

$$\int_0^{\frac{2}{3}} \int_x^{1-\frac{1}{2}x} \int_0^{2-x-2y} dz dy dx + \int_0^{\frac{2}{3}} \int_y^{1-\frac{1}{2}y} \int_0^{2-2x-y} dz dx dy = \frac{4}{9}$$

23. Find the volume of the bounded region determined by  $x \geq 0, y \geq 0, z \geq 0$ , and  $\frac{1}{7}x + y + \frac{1}{3}z = 1$ , and  $x + \frac{1}{7}y + \frac{1}{3}z = 1$ .

Answer:

$$\int_0^{\frac{7}{8}} \int_x^{1-\frac{1}{7}x} \int_0^{3-\frac{3}{7}x-3y} dz dy dx + \int_0^{\frac{7}{8}} \int_y^{1-\frac{1}{7}y} \int_0^{3-3x-\frac{3}{7}y} dz dx dy = \frac{7}{8}$$

24. Find the mass of the solid determined by  $25x^2 + 4y^2 \leq 9, z \geq 0$ , and  $z = x + 2$  if the density is  $\rho(x, y, z) = x$ .

Answer:

$$\int_{-\frac{3}{5}}^{\frac{3}{5}} \int_{-\frac{1}{2}\sqrt{(9-25x^2)}}^{\frac{1}{2}\sqrt{(9-25x^2)}} \int_0^{x+2} (x) dz dy dx = \frac{81}{1000}\pi$$

25. Find  $\int_0^1 \int_0^{35-5z} \int_{\frac{1}{5}x}^{7-z} (7-z) \cos(y^2) dy dx dz$ .

Answer:

You need to interchange the order of integration.  $\int_0^1 \int_0^{7-z} \int_0^{5y} (7-z) \cos(y^2) dx dy dz = \frac{5}{4} \cos 36 - \frac{5}{4} \cos 49$

26. Find  $\int_0^2 \int_0^{12-3z} \int_{\frac{1}{3}x}^{4-z} (4-z) \exp(y^2) dy dx dz$ .

Answer:

You need to interchange the order of integration.  $\int_0^2 \int_0^{4-z} \int_0^{3y} (4-z) \exp(y^2) dx dy dz = -\frac{3}{4}e^4 - 9 + \frac{3}{4}e^{16}$

27. Find  $\int_0^2 \int_0^{25-5z} \int_{\frac{1}{5}y}^{5-z} (5-z) \exp(x^2) dx dy dz$ .

Answer:

You need to interchange the order of integration.

$$\int_0^2 \int_0^{5-z} \int_0^{5x} (5-z) \exp(x^2) dy dx dz = -\frac{5}{4}e^9 - 20 + \frac{5}{4}e^{25}$$

28. Find  $\int_0^1 \int_0^{10-2z} \int_{\frac{1}{2}y}^{5-z} \frac{\sin x}{x} dx dy dz$ .

Answer:

You need to interchange the order of integration.

$$\int_0^1 \int_0^{5-z} \int_0^{2x} \frac{\sin x}{x} dy dx dz =$$

$$-2 \sin 1 \cos 5 + 2 \cos 1 \sin 5 + 2 - 2 \sin 5$$

29. Find  $\int_0^{20} \int_0^2 \int_{\frac{1}{5}y}^{6-z} \frac{\sin x}{x} dx dz dy + \int_{20}^{30} \int_0^{\frac{6-\frac{1}{5}y}} \int_{\frac{1}{5}y}^{6-z} \frac{\sin x}{x} dx dz dy$ .

Answer:

You need to interchange the order of integration.

$$\begin{aligned} \int_0^2 \int_0^{30-5z} \int_{\frac{1}{5}y}^{6-z} \frac{\sin x}{x} dx dy dz &= \int_0^2 \int_0^{6-z} \int_0^{5x} \frac{\sin x}{x} dy dx dz \\ &= -5 \sin 2 \cos 6 + 5 \cos 2 \sin 6 + 10 - 5 \sin 6 \end{aligned}$$



## 21.5 Exercises

1. Evaluate the integral  $\int_2^4 \int_2^{2x} \int_{2y}^x dz dy dx$
2. Find  $\int_0^3 \int_0^{2-5x} \int_0^{2-x-2y} 2x dz dy dx$
3. Find  $\int_0^2 \int_0^{1-3x} \int_0^{3-3x-2y} x dz dy dx$
4. Evaluate the integral  $\int_2^5 \int_4^{3x} \int_{4y}^x (x-y) dz dy dx$
5. Evaluate the integral  $\int_0^\pi \int_0^{3y} \int_0^{y+z} \cos(x+y) dx dz dy$
6. Evaluate the integral  $\int_0^\pi \int_0^{4y} \int_0^{y+z} \sin(x+y) dx dz dy$
7. Fill in the missing limits.  $\int_0^1 \int_0^z \int_0^z f(x, y, z) dx dy dz = \int_?^? \int_?^? \int_?^? f(x, y, z) dx dz dy$ ,  
 $\int_0^1 \int_0^z \int_0^{2z} f(x, y, z) dx dy dz = \int_?^? \int_?^? \int_?^? f(x, y, z) dy dz dx$ ,  
 $\int_0^1 \int_0^z \int_0^z f(x, y, z) dx dy dz = \int_?^? \int_?^? \int_?^? f(x, y, z) dz dy dx$ ,  
 $\int_0^1 \int_{z/2}^{\sqrt{z}} \int_0^{y+z} f(x, y, z) dx dy dz = \int_?^? \int_?^? \int_?^? f(x, y, z) dx dz dy$ ,  
 $\int_4^6 \int_2^6 \int_0^4 f(x, y, z) dx dy dz = \int_?^? \int_?^? \int_?^? f(x, y, z) dz dy dx$ .
8. Find the volume of  $R$  where  $R$  is the bounded region formed by the plane  $\frac{1}{5}x + \frac{1}{3}y + \frac{1}{4}z = 1$  and the planes  $x = 0, y = 0, z = 0$ .
9. Find the volume of  $R$  where  $R$  is the bounded region formed by the plane  $\frac{1}{4}x + \frac{1}{2}y + \frac{1}{4}z = 1$  and the planes  $x = 0, y = 0, z = 0$ .
10. Find the mass of the bounded region,  $R$  formed by the plane  $\frac{1}{4}x + \frac{1}{3}y + \frac{1}{2}z = 1$  and the planes  $x = 0, y = 0, z = 0$  if the density is  $\rho(x, y, z) = y + z$
11. Find the mass of the bounded region,  $R$  formed by the plane  $\frac{1}{4}x + \frac{1}{2}y + \frac{1}{5}z = 1$  and the planes  $x = 0, y = 0, z = 0$  if the density is  $\rho(x, y, z) = y$
12. Here is an iterated integral:  $\int_0^2 \int_0^{1-\frac{1}{2}x} \int_0^{x^2} dz dy dx$ . Write as an iterated integral in the following orders:  $dz dx dy$ ,  $dx dz dy$ ,  $dx dy dz$ ,  $dy dx dz$ ,  $dy dz dx$ .
13. Find the volume of the bounded region determined by  $2y + z = 3, x = 9 - y^2, y = 0, x = 0$ .
14. Find the volume of the bounded region determined by  $3y + 2z = 5, x = 9 - y^2, y = 0, x = 0$ .  
Your answer should be  $\frac{11\ 525}{648}$
15. Find the volume of the bounded region determined by  $5y + 2z = 3, x = 9 - y^2, y = 0, x = 0$ .
16. Find the volume of the region bounded by  $x^2 + y^2 = 25, z = x, z = 0$ , and  $x \geq 0$ .  
Your answer should be  $\frac{250}{3}$ .
17. Find the volume of the region bounded by  $x^2 + y^2 = 9, z = 3x, z = 0$ , and  $x \geq 0$ .
18. Find the volume of the region determined by the intersection of the two cylinders,  $x^2 + y^2 \leq 16$  and  $y^2 + z^2 \leq 16$ .

19. Find the total mass of the bounded solid determined by  $z = 4 - x^2 - y^2$  and  $x, y, z \geq 0$  if the mass is given by  $\rho(x, y, z) = y$
20. Find the total mass of the bounded solid determined by  $z = 9 - x^2 - y^2$  and  $x, y, z \geq 0$  if the mass is given by  $\rho(x, y, z) = z^2$
21. Find the volume of the region bounded by  $x^2 + y^2 = 4, z = 0, z = 5 - y$
22. Find the volume of the bounded region determined by  $x \geq 0, y \geq 0, z \geq 0$ , and  $\frac{1}{7}x + \frac{1}{3}y + \frac{1}{3}z = 1$ , and  $\frac{1}{3}x + \frac{1}{7}y + \frac{1}{3}z = 1$ .
23. Find the volume of the bounded region determined by  $x \geq 0, y \geq 0, z \geq 0$ , and  $\frac{1}{5}x + \frac{1}{3}y + z = 1$ , and  $\frac{1}{3}x + \frac{1}{5}y + z = 1$ .
24. Find the mass of the solid determined by  $16x^2 + 4y^2 \leq 9, z \geq 0$ , and  $z = x + 2$  if the density is  $\rho(x, y, z) = z$ .
25. Find  $\int_0^2 \int_0^{6-2z} \int_{\frac{1}{2}x}^{3-z} (3-z) \cos(y^2) dy dx dz$ .
26. Find  $\int_0^1 \int_0^{18-3z} \int_{\frac{1}{3}x}^{6-z} (6-z) \exp(y^2) dy dx dz$ .
27. Find  $\int_0^2 \int_0^{24-4z} \int_{\frac{1}{4}y}^{6-z} (6-z) \exp(x^2) dx dy dz$ .
28. Find  $\int_0^1 \int_0^{12-4z} \int_{\frac{1}{4}y}^{3-z} \frac{\sin x}{x} dx dy dz$ .
29. Find  $\int_0^{20} \int_0^1 \int_{\frac{1}{5}y}^{5-z} \frac{\sin x}{x} dx dz dy + \int_{20}^{25} \int_0^{5-\frac{1}{5}y} \int_{\frac{1}{5}y}^{5-z} \frac{\sin x}{x} dx dz dy$ . **Hint:** You might try doing it in the order,  $dy dx dz$

## 21.6 Different Coordinates

As mentioned above, the fundamental concept of an integral is a sum of things of the form  $f(\mathbf{x}) dV$  where  $dV$  is an “infinitesimal” chunk of volume located at the point,  $\mathbf{x}$ . Up to now, this infinitesimal chunk of volume has had the form of a box with sides  $dx_1, \dots, dx_n$  so  $dV = dx_1 dx_2 \cdots dx_n$  but its form is not important. It could just as well be an infinitesimal parallelepiped for example. In what follows, this is what it will be.

First recall the following fundamental definition on Page 460.

**Definition 21.6.1** Let  $\mathbf{u}_1, \dots, \mathbf{u}_p$  be vectors in  $\mathbb{R}^k$ . The parallelepiped determined by these vectors will be denoted by  $P(\mathbf{u}_1, \dots, \mathbf{u}_p)$  and it is defined as

$$P(\mathbf{u}_1, \dots, \mathbf{u}_p) \equiv \left\{ \sum_{j=1}^p s_j \mathbf{u}_j : s_j \in [0, 1] \right\}.$$

Now define the volume of this parallelepiped.

$$\text{volume of } P(\mathbf{u}_1, \dots, \mathbf{u}_p) \equiv (\det(\mathbf{u}_i \cdot \mathbf{u}_j))^{1/2}.$$

The dot product is used to determine this volume of a parallelepiped spanned by the given vectors and you should note that it is only the dot product that matters. Now consider spherical coordinates,  $\rho, \phi$ , and  $\theta$ . Recall there is a relationship between these coordinates and rectangular coordinates given by

$$x = \rho \sin \phi \cos \theta, y = \rho \sin \phi \sin \theta, z = \rho \cos \phi \quad (21.8)$$

where  $\phi \in [0, \pi]$ ,  $\theta \in [0, 2\pi)$ , and  $\rho > 0$ . Thus  $(\rho, \phi, \theta)$  is a point in  $\mathbb{R}^3$ , more specifically in the set

$$U = (0, \infty) \times [0, \pi] \times [0, 2\pi)$$

and corresponding to such a  $(\rho, \phi, \theta) \in U$  there exists a unique point,  $(x, y, z) \in V$  where  $V$  consists of all points of  $\mathbb{R}^3$  other than the origin,  $(0, 0, 0)$ . This  $(x, y, z)$  determines a unique point in three dimensional space as mentioned earlier. Suppose at the point  $(\rho_0, \phi_0, \theta_0) \in U$ , there is an infinitesimal box having sides  $d\rho, d\phi, d\theta$ . Then this little box would correspond to something in  $V$ . What? Consider the mapping from  $U$  to  $V$  defined by

$$\mathbf{x} = \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} \rho \sin \phi \cos \theta \\ \rho \sin \phi \sin \theta \\ \rho \cos \phi \end{pmatrix} = \mathbf{f}(\rho, \phi, \theta) \quad (21.9)$$

which takes a point,  $(\rho, \phi, \theta)$  in  $U$  and sends it to the point in  $V$  which is identified as  $(x, y, z)^T \equiv \mathbf{x}$ . What happens to an element of the infinitesimal box, located at  $(\rho_0, \phi_0, \theta_0)$ ? Such an element is of the form

$$(\rho_0 + s_1 d\rho, \phi_0 + s_2 d\phi, \theta_0 + s_3 d\theta),$$

where  $s_i \geq 0$  and  $\sum_i s_i \leq 1$ . Also, from the definition of the derivative,

$$\begin{aligned} \mathbf{f}(\rho_0 + s_1 d\rho, \phi_0 + s_2 d\phi, \theta_0 + s_3 d\theta) - \mathbf{f}(\rho_0, \phi_0, \theta_0) = \\ D\mathbf{f}(\rho_0, \phi_0, \theta_0) \begin{pmatrix} s_1 d\rho \\ s_2 d\phi \\ s_3 d\theta \end{pmatrix} + \mathbf{o} \begin{pmatrix} s_1 d\rho \\ s_2 d\phi \\ s_3 d\theta \end{pmatrix} \end{aligned}$$

where the last term may be taken equal to  $\mathbf{0}$  because the vector,  $(s_1 d\rho, s_2 d\phi, s_3 d\theta)^T$  is infinitesimal meaning nothing precise but conveying the idea that it is surpassingly small. Therefore, an element of this infinitesimal box is sent to the vector,

$$\begin{aligned} & \overbrace{\left( \frac{\partial \mathbf{x}(\rho_0, \phi_0, \theta_0)}{\partial \rho}, \frac{\partial \mathbf{x}(\rho_0, \phi_0, \theta_0)}{\partial \phi}, \frac{\partial \mathbf{x}(\rho_0, \phi_0, \theta_0)}{\partial \theta} \right)}^{= D\mathbf{f}(\rho_0, \phi_0, \theta_0)} \begin{pmatrix} s_1 d\rho \\ s_2 d\phi \\ s_3 d\theta \end{pmatrix} = \\ & s_1 \frac{\partial \mathbf{x}(\rho_0, \phi_0, \theta_0)}{\partial \rho} d\rho + s_2 \frac{\partial \mathbf{x}(\rho_0, \phi_0, \theta_0)}{\partial \phi} d\phi + s_3 \frac{\partial \mathbf{x}(\rho_0, \phi_0, \theta_0)}{\partial \theta} d\theta \end{aligned}$$

an element of the infinitesimal parallelepiped determined by the vectors

$$\left\{ \frac{\partial \mathbf{x}(\rho_0, \phi_0, \theta_0)}{\partial \rho} d\rho, \frac{\partial \mathbf{x}(\rho_0, \phi_0, \theta_0)}{\partial \phi} d\phi, \frac{\partial \mathbf{x}(\rho_0, \phi_0, \theta_0)}{\partial \theta} d\theta \right\}.$$

The situation is no different for general coordinate systems. In general,  $\mathbf{x} = \mathbf{f}(\mathbf{u})$  where  $\mathbf{u} \in U$ , a subset of  $\mathbb{R}^n$  and  $\mathbf{x}$  is a point in  $V$ , a subset of  $n$  dimensional space. Thus, letting the Cartesian coordinates of  $\mathbf{x}$  be given by  $\mathbf{x} = (x_1, \dots, x_n)^T$ , each  $x_i$  being a function of  $\mathbf{u}$ , an infinitesimal box located at  $\mathbf{u}_0$  corresponds to an infinitesimal parallelepiped located at  $\mathbf{f}(\mathbf{u}_0)$  which is determined by the  $n$  vectors  $\left\{ \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_i} du_i \right\}_{i=1}^n$ . From Definition 21.6.1, the volume of this infinitesimal parallelepiped located at  $\mathbf{f}(\mathbf{u}_0)$  is given by

$$\det \left( \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_i} du_i \cdot \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_j} du_j \right)^{1/2} \quad (21.10)$$

in which there is no sum on the repeated index. Now in general if there are  $n$  vectors in  $\mathbb{R}^n$ ,  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ ,

$$\det(\mathbf{v}_i \cdot \mathbf{v}_j)^{1/2} = |\det(\mathbf{v}_1, \dots, \mathbf{v}_n)| \quad (21.11)$$

where this last matrix is the  $n \times n$  matrix which has the  $i^{th}$  column equal to  $\mathbf{v}_i$ . The reason for this is that the matrix whose  $ij^{th}$  entry is  $\mathbf{v}_i \cdot \mathbf{v}_j$  is just the product of the two matrices,

$$\begin{pmatrix} \mathbf{v}_1^T \\ \vdots \\ \mathbf{v}_n^T \end{pmatrix} (\mathbf{v}_1, \dots, \mathbf{v}_n)$$

where the first on the left is the matrix having the  $i^{th}$  row equal to  $\mathbf{v}_i^T$  while the matrix on the right is just the matrix having the  $i^{th}$  column equal to  $\mathbf{v}_i$ . Therefore, since the determinant of a matrix equals the determinant of its transpose,

$$\begin{aligned} \det(\mathbf{v}_i \cdot \mathbf{v}_j) &= \det \left( \begin{pmatrix} \mathbf{v}_1^T \\ \vdots \\ \mathbf{v}_n^T \end{pmatrix} (\mathbf{v}_1, \dots, \mathbf{v}_n) \right) \\ &= \det(\mathbf{v}_1, \dots, \mathbf{v}_n)^2 \end{aligned}$$

and so taking square roots yields (21.11). Therefore, from the properties of determinants, (21.10) equals

$$\begin{aligned} &\left| \det \left( \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_1} du_1, \dots, \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_n} du_n \right) \right| = \\ &\left| \det \left( \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_1}, \dots, \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_n} \right) \right| du_1 \cdots du_n \end{aligned}$$

and this is the infinitesimal chunk of volume corresponding to the point  $\mathbf{f}(\mathbf{u}_0)$  in  $V$ .

**Definition 21.6.2** Let  $\mathbf{x} = \mathbf{f}(\mathbf{u})$  be as described above. Then the symbol,  $\frac{\partial(x_1, \dots, x_n)}{\partial(u_1, \dots, u_n)}$ , called the Jacobian determinant, is defined by

$$\left| \det \left( \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_1}, \dots, \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_n} \right) \right| \equiv \frac{\partial(x_1, \dots, x_n)}{\partial(u_1, \dots, u_n)}.$$

Also, the symbol,  $\left| \frac{\partial(x_1, \dots, x_n)}{\partial(u_1, \dots, u_n)} \right| du_1 \cdots du_n$  is called the volume element.

This has given motivation for the following fundamental procedure often called the change of variables formula.

**Procedure 21.6.3** Suppose  $U$  is a subset of  $\mathbb{R}^n$  and suppose  $\mathbf{f} : U \rightarrow V$  is a  $C^1$  function which is one to one.<sup>2</sup> Then if  $h : V \rightarrow \mathbb{R}$ ,

$$\int_U h(\mathbf{f}(\mathbf{u})) \frac{\partial(x_1, \dots, x_n)}{\partial(u_1, \dots, u_n)} dV = \int_V h(\mathbf{x}) dV.$$

<sup>2</sup>This will cause non overlapping infinitesimal boxes in  $U$  to be mapped to non overlapping infinitesimal parallelepipeds in  $V$ .

Also, in the context of the Riemann integral we should say more about the sets,  $U$  and  $V$  in any case the function,  $h$ . These conditions are mainly technical however, and since a mathematically respectable treatment will not be attempted for this theorem, I think it best to give a memorable version of it which is essentially correct in all examples of interest.

Now return to Spherical coordinates. In this case, it is necessary to find the absolute value of

$$\det \left( \frac{\partial \mathbf{x}(\rho_0, \phi_0, \theta_0)}{\partial \rho}, \frac{\partial \mathbf{x}(\rho_0, \phi_0, \theta_0)}{\partial \phi}, \frac{\partial \mathbf{x}(\rho_0, \phi_0, \theta_0)}{\partial \theta} \right)$$

which equals

$$\det \begin{pmatrix} \sin \phi \cos \theta & \rho \cos \phi \cos \theta & -\rho \sin \phi \sin \theta \\ \sin \phi \sin \theta & \rho \cos \phi \sin \theta & \rho \sin \phi \cos \theta \\ \cos \phi & -\rho \sin \phi & 0 \end{pmatrix} = \rho^2 \sin \phi$$

which is positive because  $\phi \in [0, \pi]$ .

**Example 21.6.4** Find the volume of a ball,  $B_R$  of radius  $R$ .

In this case,  $U = (0, R] \times [0, \pi] \times [0, 2\pi)$  and use spherical coordinates. Then (21.9) yields a set in  $\mathbb{R}^3$  which clearly differs from the ball of radius  $R$  only by a set having volume equal to zero. It leaves out the point at the origin is all. Therefore, the volume of the ball is

$$\begin{aligned} \int_{B_R} 1 dV &= \int_U \rho^2 \sin \phi dV \\ &= \int_0^R \int_0^\pi \int_0^{2\pi} \rho^2 \sin \phi d\theta d\phi d\rho = \frac{4}{3} R^3 \pi. \end{aligned}$$

The reason this was effortless, is that the ball,  $B_R$  is realized as a box in terms of the spherical coordinates. Remember what was pointed out earlier about setting up iterated integrals over boxes.

**Example 21.6.5** Find the volume element for cylindrical coordinates.

In cylindrical coordinates,

$$\begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} r \cos \theta \\ r \sin \theta \\ z \end{pmatrix}$$

Therefore, the Jacobian determinant is

$$\det \begin{pmatrix} \cos \theta & -r \sin \theta & 0 \\ \sin \theta & r \cos \theta & 0 \\ 0 & 0 & 1 \end{pmatrix} = r.$$

It follows the volume element in cylindrical coordinates is  $r d\theta dr dz$ .

**Example 21.6.6** This example uses spherical coordinates to verify an important conclusion about gravitational force. Let the hollow sphere,  $H$  be defined by  $a^2 \leq x^2 + y^2 + z^2 \leq b^2$  and suppose this hollow sphere has constant density taken to equal 1. Now place a unit mass at the point  $(0, 0, z_0)$  where  $|z_0| \in [a, b]$ . Show the force of gravity acting on this unit mass is  $\left( \alpha G \int \int \int_H \frac{(z - z_0)}{[x^2 + y^2 + (z - z_0)^2]^{3/2}} dV \right) \mathbf{k}$  and then show that if  $|z_0| > b$  then the force of gravity acting on this point mass is the same as if the entire mass of the hollow sphere were placed at the origin, while if  $|z_0| < a$ , the total force acting on the point mass from gravity equals zero. Here  $G$  is the gravitation constant and  $\alpha$  is the density. In particular, this shows that the force a planet exerts on an object is as though the entire mass of the planet were situated at its center<sup>3</sup>.

<sup>3</sup>This was shown by Newton in 1685 and allowed him to assert his law of gravitation applied to the planets as though they were point masses. It was a major accomplishment.

Without loss of generality, assume  $z_0 > 0$ . Let  $dV$  be a little chunk of material located at the point  $(x, y, z)$  of  $H$  the hollow sphere. Then according to Newton's law of gravity, the force this small chunk of material exerts on the given point mass equals

$$\frac{x\mathbf{i} + y\mathbf{j} + (z - z_0)\mathbf{k}}{|x\mathbf{i} + y\mathbf{j} + (z - z_0)\mathbf{k}|} \frac{1}{\left(x^2 + y^2 + (z - z_0)^2\right)} G\alpha dV =$$

$$(x\mathbf{i} + y\mathbf{j} + (z - z_0)\mathbf{k}) \frac{1}{\left(x^2 + y^2 + (z - z_0)^2\right)^{3/2}} G\alpha dV$$

Therefore, the total force is

$$\int \int \int_H (x\mathbf{i} + y\mathbf{j} + (z - z_0)\mathbf{k}) \frac{1}{\left(x^2 + y^2 + (z - z_0)^2\right)^{3/2}} G\alpha dV.$$

By the symmetry of the sphere, the  $\mathbf{i}$  and  $\mathbf{j}$  components will cancel out when the integral is taken. This is because there is the same amount of stuff for negative  $x$  and  $y$  as there is for positive  $x$  and  $y$ . Hence what remains is

$$\alpha G \mathbf{k} \int \int \int_H \frac{(z - z_0)}{\left[x^2 + y^2 + (z - z_0)^2\right]^{3/2}} dV$$

as claimed. Now for the interesting part, the integral is evaluated. In spherical coordinates this integral is.

$$\int_0^{2\pi} \int_a^b \int_0^\pi \frac{(\rho \cos \phi - z_0) \rho^2 \sin \phi}{(\rho^2 + z_0^2 - 2\rho z_0 \cos \phi)^{3/2}} d\phi d\rho d\theta. \quad (21.12)$$

Rewrite the inside integral and use integration by parts to obtain this inside integral equals

$$\frac{1}{2z_0} \int_0^\pi (\rho^2 \cos \phi - \rho z_0) \frac{(2z_0 \rho \sin \phi)}{(\rho^2 + z_0^2 - 2\rho z_0 \cos \phi)^{3/2}} d\phi =$$

$$\frac{1}{2z_0} \left( -2 \frac{-\rho^2 - \rho z_0}{\sqrt{(\rho^2 + z_0^2 + 2\rho z_0)}} + 2 \frac{\rho^2 - \rho z_0}{\sqrt{(\rho^2 + z_0^2 - 2\rho z_0)}} - \int_0^\pi 2\rho^2 \frac{\sin \phi}{\sqrt{(\rho^2 + z_0^2 - 2\rho z_0 \cos \phi)}} d\phi \right). \quad (21.13)$$

There are some cases to consider here.

First suppose  $z_0 < a$  so the point is on the inside of the hollow sphere and it is always the case that  $\rho > z_0$ . Then in this case, the two first terms reduce to

$$\frac{2\rho(\rho + z_0)}{\sqrt{(\rho + z_0)^2}} + \frac{2\rho(\rho - z_0)}{\sqrt{(\rho - z_0)^2}} = \frac{2\rho(\rho + z_0)}{(\rho + z_0)} + \frac{2\rho(\rho - z_0)}{\rho - z_0} = 4\rho$$

and so the expression in (21.13) equals

$$\frac{1}{2z_0} \left( 4\rho - \int_0^\pi 2\rho^2 \frac{\sin \phi}{\sqrt{(\rho^2 + z_0^2 - 2\rho z_0 \cos \phi)}} d\phi \right)$$

$$= \frac{1}{2z_0} \left( 4\rho - \frac{1}{z_0} \int_0^\pi \rho \frac{2\rho z_0 \sin \phi}{\sqrt{(\rho^2 + z_0^2 - 2\rho z_0 \cos \phi)}} d\phi \right)$$

$$\begin{aligned}
&= \frac{1}{2z_0} \left( 4\rho - \frac{2\rho}{z_0} (\rho^2 + z_0^2 - 2\rho z_0 \cos \phi)^{1/2} \Big|_0^\pi \right) \\
&= \frac{1}{2z_0} \left( 4\rho - \frac{2\rho}{z_0} [(\rho + z_0) - (\rho - z_0)] \right) = 0.
\end{aligned}$$

Therefore, in this case the inner integral of (21.12) equals zero and so the original integral will also be zero.

The other case is when  $z_0 > b$  and so it is always the case that  $z_0 > \rho$ . In this case the first two terms of (21.13) are

$$\frac{2\rho(\rho + z_0)}{\sqrt{(\rho + z_0)^2}} + \frac{2\rho(\rho - z_0)}{\sqrt{(\rho - z_0)^2}} = \frac{2\rho(\rho + z_0)}{(\rho + z_0)} + \frac{2\rho(\rho - z_0)}{z_0 - \rho} = 0.$$

Therefore in this case, (21.13) equals

$$\begin{aligned}
&\frac{1}{2z_0} \left( - \int_0^\pi 2\rho^2 \frac{\sin \phi}{\sqrt{(\rho^2 + z_0^2 - 2\rho z_0 \cos \phi)}} d\phi \right) \\
&= \frac{-\rho}{2z_0^2} \left( \int_0^\pi \frac{2\rho z_0 \sin \phi}{\sqrt{(\rho^2 + z_0^2 - 2\rho z_0 \cos \phi)}} d\phi \right)
\end{aligned}$$

which equals

$$\begin{aligned}
&\frac{-\rho}{z_0^2} \left( (\rho^2 + z_0^2 - 2\rho z_0 \cos \phi)^{1/2} \Big|_0^\pi \right) \\
&= \frac{-\rho}{z_0^2} [(\rho + z_0) - (z_0 - \rho)] = -\frac{2\rho^2}{z_0^2}.
\end{aligned}$$

Thus the inner integral of (21.12) reduces to the above simple expression. Therefore, (21.12) equals

$$\int_0^{2\pi} \int_a^b \left( -\frac{2}{z_0^2} \rho^2 \right) d\rho d\theta = -\frac{4}{3} \pi \frac{b^3 - a^3}{z_0^2}$$

and so

$$\alpha G \mathbf{k} \int \int \int_H \frac{(z - z_0)}{[x^2 + y^2 + (z - z_0)^2]^{3/2}} dV = \alpha G \mathbf{k} \left( -\frac{4}{3} \pi \frac{b^3 - a^3}{z_0^2} \right) = -\mathbf{k} G \frac{\text{total mass}}{z_0^2}.$$

## 21.7 Exercises With Answers

1. Find the area of the bounded region,  $R$ , determined by  $3x+3y=1$ ,  $3x+3y=8$ ,  $y=3x$ , and  $y=4x$ .

Answer:

Let  $u = \frac{y}{x}$ ,  $v = 3x + 3y$ . Then solving these equations for  $x$  and  $y$  yields

$$\left\{ x = \frac{1}{3} \frac{v}{1+u}, y = \frac{1}{3} u \frac{v}{1+u} \right\}.$$

Now

$$\frac{\partial(x, y)}{\partial(u, v)} = \det \begin{pmatrix} -\frac{1}{3} \frac{v}{(1+u)^2} & \frac{1}{3+3u} \\ \frac{1}{3} \frac{v}{(1+u)^2} & \frac{1}{3} \frac{u}{1+u} \end{pmatrix} = -\frac{1}{9} \frac{v}{(1+u)^2}.$$

Also,  $u \in [3, 4]$  while  $v \in [1, 8]$ . Therefore,

$$\begin{aligned}\int_R dV &= \int_3^4 \int_1^8 \left| -\frac{1}{9} \frac{v}{(1+u)^2} \right| dv du = \\ &= \int_3^4 \int_1^8 \frac{1}{9} \frac{v}{(1+u)^2} dv du = \frac{7}{40}\end{aligned}$$

2. Find the area of the bounded region,  $R$ , determined by  $5x + y = 1$ ,  $5x + y = 9$ ,  $y = 2x$ , and  $y = 5x$ .

Answer:

Let  $u = \frac{y}{x}$ ,  $v = 5x + y$ . Then solving these equations for  $x$  and  $y$  yields

$$\left\{ x = \frac{v}{5+u}, y = u \frac{v}{5+u} \right\}.$$

Now

$$\frac{\partial(x, y)}{\partial(u, v)} = \det \begin{pmatrix} -\frac{v}{(5+u)^2} & \frac{1}{5+u} \\ 5\frac{v}{(5+u)^2} & \frac{u}{5+u} \end{pmatrix} = -\frac{v}{(5+u)^2}.$$

Also,  $u \in [2, 5]$  while  $v \in [1, 9]$ . Therefore,

$$\int_R dV = \int_2^5 \int_1^9 \left| -\frac{v}{(5+u)^2} \right| dv du = \int_2^5 \int_1^9 \frac{v}{(5+u)^2} dv du = \frac{12}{7}$$

3. A solid,  $R$  is determined by  $5x + 3y = 4$ ,  $5x + 3y = 9$ ,  $y = 2x$ , and  $y = 5x$  and the density is  $\rho = x$ . Find the total mass of  $R$ .

Answer:

Let  $u = \frac{y}{x}$ ,  $v = 5x + 3y$ . Then solving these equations for  $x$  and  $y$  yields

$$\left\{ x = \frac{v}{5+3u}, y = u \frac{v}{5+3u} \right\}.$$

Now

$$\frac{\partial(x, y)}{\partial(u, v)} = \det \begin{pmatrix} -3\frac{v}{(5+3u)^2} & \frac{1}{5+3u} \\ 5\frac{v}{(5+3u)^2} & \frac{u}{5+3u} \end{pmatrix} = -\frac{v}{(5+3u)^2}.$$

Also,  $u \in [2, 5]$  while  $v \in [4, 9]$ . Therefore,

$$\begin{aligned}\int_R \rho dV &= \int_2^5 \int_4^9 \frac{v}{5+3u} \left| -\frac{v}{(5+3u)^2} \right| dv du = \\ &= \int_2^5 \int_4^9 \left( \frac{v}{5+3u} \right) \left( \frac{v}{(5+3u)^2} \right) dv du = \frac{4123}{19360}.\end{aligned}$$



4. A solid,  $R$  is determined by  $2x + 2y = 1$ ,  $2x + 2y = 10$ ,  $y = 4x$ , and  $y = 5x$  and the density is  $\rho = x + 1$ . Find the total mass of  $R$ .

Answer:

Let  $u = \frac{y}{x}$ ,  $v = 2x + 2y$ . Then solving these equations for  $x$  and  $y$  yields

$$\left\{ x = \frac{1}{2} \frac{v}{1+u}, y = \frac{1}{2} u \frac{v}{1+u} \right\}.$$

Now

$$\frac{\partial(x, y)}{\partial(u, v)} = \det \begin{pmatrix} -\frac{1}{2} \frac{v}{(1+u)^2} & \frac{1}{2+2u} \\ \frac{1}{2} \frac{v}{(1+u)^2} & \frac{1}{2} \frac{u}{1+u} \end{pmatrix} = -\frac{1}{4} \frac{v}{(1+u)^2}.$$

Also,  $u \in [4, 5]$  while  $v \in [1, 10]$ . Therefore,

$$\begin{aligned} \int_R \rho dV &= \int_4^5 \int_1^{10} (x+1) \left| -\frac{1}{4} \frac{v}{(1+u)^2} \right| dv du \\ &= \int_4^5 \int_1^{10} (x+1) \left( \frac{1}{4} \frac{v}{(1+u)^2} \right) dv du \end{aligned}$$

5. A solid,  $R$  is determined by  $4x + 2y = 1$ ,  $4x + 2y = 9$ ,  $y = x$ , and  $y = 6x$  and the density is  $\rho = y^{-1}$ . Find the total mass of  $R$ .

Answer:

Let  $u = \frac{y}{x}$ ,  $v = 4x + 2y$ . Then solving these equations for  $x$  and  $y$  yields

$$\left\{ x = \frac{1}{2} \frac{v}{2+u}, y = \frac{1}{2} u \frac{v}{2+u} \right\}.$$

Now

$$\frac{\partial(x, y)}{\partial(u, v)} = \det \begin{pmatrix} -\frac{1}{2} \frac{v}{(2+u)^2} & \frac{1}{4+2u} \\ \frac{v}{(2+u)^2} & \frac{1}{2} \frac{u}{2+u} \end{pmatrix} = -\frac{1}{4} \frac{v}{(2+u)^2}.$$

Also,  $u \in [1, 6]$  while  $v \in [1, 9]$ . Therefore,

$$\int_R \rho dV = \int_1^6 \int_1^9 \left( \frac{1}{2} u \frac{v}{2+u} \right)^{-1} \left| -\frac{1}{4} \frac{v}{(2+u)^2} \right| dv du = -4 \ln 2 + 4 \ln 3$$

6. Find the volume of the region,  $E$ , bounded by the ellipsoid,  $\frac{1}{4}x^2 + \frac{1}{9}y^2 + \frac{1}{49}z^2 = 1$ .

Answer:

Let  $u = \frac{1}{2}x$ ,  $v = \frac{1}{3}y$ ,  $w = \frac{1}{7}z$ . Then  $(u, v, w)$  is a point in the unit ball,  $B$ . Therefore,

$$\int_B \frac{\partial(x, y, z)}{\partial(u, v, w)} dV = \int_E dV.$$

But  $\frac{\partial(x, y, z)}{\partial(u, v, w)} = 42$  and so the answer is

$$(\text{volume of } B) \times 42 = \frac{4}{3}\pi 42 = 56\pi.$$

7. Here are three vectors.  $(4, 1, 4)^T$ ,  $(5, 0, 4)^T$ , and  $(3, 1, 5)^T$ . These vectors determine a parallelepiped,  $R$ , which is occupied by a solid having density  $\rho = x$ . Find the mass of this solid.

Answer:

Let  $\begin{pmatrix} 4 & 5 & 3 \\ 1 & 0 & 1 \\ 4 & 4 & 5 \end{pmatrix} \begin{pmatrix} u \\ v \\ w \end{pmatrix} = \begin{pmatrix} x \\ y \\ z \end{pmatrix}$ . Then this maps the unit cube,

$$Q \equiv [0, 1] \times [0, 1] \times [0, 1]$$

onto  $R$  and

$$\frac{\partial(x, y, z)}{\partial(u, v, w)} = \left| \det \begin{pmatrix} 4 & 5 & 3 \\ 1 & 0 & 1 \\ 4 & 4 & 5 \end{pmatrix} \right| = |-9| = 9$$

so the mass is

$$\begin{aligned} \int_R x \, dV &= \int_Q (4u + 5v + 3w) (9) \, dV \\ &= \int_0^1 \int_0^1 \int_0^1 (4u + 5v + 3w) (9) \, du \, dv \, dw = 54 \end{aligned}$$

8. Here are three vectors.  $(3, 2, 6)^T$ ,  $(4, 1, 6)^T$ , and  $(2, 2, 7)^T$ . These vectors determine a parallelepiped,  $R$ , which is occupied by a solid having density  $\rho = y$ . Find the mass of this solid.

Answer:

Let  $\begin{pmatrix} 3 & 4 & 2 \\ 2 & 1 & 2 \\ 6 & 6 & 7 \end{pmatrix} \begin{pmatrix} u \\ v \\ w \end{pmatrix} = \begin{pmatrix} x \\ y \\ z \end{pmatrix}$ . Then this maps the unit cube,

$$Q \equiv [0, 1] \times [0, 1] \times [0, 1]$$

onto  $R$  and

$$\frac{\partial(x, y, z)}{\partial(u, v, w)} = \left| \det \begin{pmatrix} 3 & 4 & 2 \\ 2 & 1 & 2 \\ 6 & 6 & 7 \end{pmatrix} \right| = |-11| = 11$$

and so the mass is

$$\begin{aligned} \int_R x \, dV &= \int_Q (2u + v + 2w) (11) \, dV \\ &= \int_0^1 \int_0^1 \int_0^1 (2u + v + 2w) (11) \, du \, dv \, dw = \frac{55}{2}. \end{aligned}$$

9. Here are three vectors.  $(2, 2, 4)^T$ ,  $(3, 1, 4)^T$ , and  $(1, 2, 5)^T$ . These vectors determine a parallelepiped,  $R$ , which is occupied by a solid having density  $\rho = y + x$ . Find the mass of this solid.

Answer:

Let  $\begin{pmatrix} 2 & 3 & 1 \\ 2 & 1 & 2 \\ 4 & 4 & 5 \end{pmatrix} \begin{pmatrix} u \\ v \\ w \end{pmatrix} = \begin{pmatrix} x \\ y \\ z \end{pmatrix}$ . Then this maps the unit cube,

$$Q \equiv [0, 1] \times [0, 1] \times [0, 1]$$

onto  $R$  and

$$\frac{\partial(x, y, z)}{\partial(u, v, w)} = \left| \det \begin{pmatrix} 2 & 3 & 1 \\ 2 & 1 & 2 \\ 4 & 4 & 5 \end{pmatrix} \right| = |-8| = 8$$

and so the mass is  $2u + 3v + w$

$$\begin{aligned} \int_R x \, dV &= \int_Q (4u + 4v + 3w) (8) \, dV \\ &= \int_0^1 \int_0^1 \int_0^1 (4u + 4v + 3w) (8) \, du \, dv \, dw = 44. \end{aligned}$$

10. Let  $D = \{(x, y) : x^2 + y^2 \leq 25\}$ . Find  $\int_D e^{36x^2+36y^2} \, dx \, dy$ .

Answer:

This is easy in polar coordinates.  $x = r \cos \theta$ ,  $y = r \sin \theta$ . Thus  $\frac{\partial(x, y)}{\partial(r, \theta)} = r$  and in terms of these new coordinates, the disk,  $D$ , is the rectangle,

$$R = \{(r, \theta) \in [0, 5] \times [0, 2\pi]\}.$$

Therefore,

$$\begin{aligned} \int_D e^{36x^2+36y^2} \, dV &= \int_R e^{36r^2} r \, dV = \\ &= \int_0^5 \int_0^{2\pi} e^{36r^2} r \, d\theta \, dr = \frac{1}{36} \pi (e^{900} - 1). \end{aligned}$$

Note you wouldn't get very far without changing the variables in this.

11. Let  $D = \{(x, y) : x^2 + y^2 \leq 9\}$ . Find  $\int_D \cos(36x^2 + 36y^2) \, dx \, dy$ .

Answer:

This is easy in polar coordinates.  $x = r \cos \theta$ ,  $y = r \sin \theta$ . Thus  $\frac{\partial(x, y)}{\partial(r, \theta)} = r$  and in terms of these new coordinates, the disk,  $D$ , is the rectangle,

$$R = \{(r, \theta) \in [0, 3] \times [0, 2\pi]\}.$$

Therefore,

$$\int_D \cos(36x^2 + 36y^2) \, dV = \int_R \cos(36r^2) r \, dV =$$

$$\int_0^3 \int_0^{2\pi} \cos(36r^2) r \, d\theta \, dr = \frac{1}{36} (\sin 324) \pi.$$

12. The ice cream in a sugar cone is described in spherical coordinates by  $\rho \in [0, 8]$ ,  $\phi \in [0, \frac{1}{4}\pi]$ ,  $\theta \in [0, 2\pi]$ . If the units are in centimeters, find the total volume in cubic centimeters of this ice cream.

Answer:

Remember that in spherical coordinates, the volume element is  $\rho^2 \sin \phi \, dV$  and so the total volume of this is  $\int_0^8 \int_0^{\frac{1}{4}\pi} \int_0^{2\pi} \rho^2 \sin \phi \, d\theta \, d\phi \, d\rho = -\frac{512}{3}\sqrt{2}\pi + \frac{1024}{3}\pi$ .

13. Find the volume between  $z = 5 - x^2 - y^2$  and  $z = \sqrt{(x^2 + y^2)}$ .

Answer:

Use cylindrical coordinates. In terms of these coordinates the shape is

$$h - r^2 \geq z \geq r, r \in \left[0, \frac{1}{2}\sqrt{21} - \frac{1}{2}\right], \theta \in [0, 2\pi].$$

Also,  $\frac{\partial(x,y,z)}{\partial(r,\theta,z)} = r$ . Therefore, the volume is

$$\int_0^{2\pi} \int_0^{\frac{1}{2}\sqrt{21}-\frac{1}{2}} \int_0^{5-r^2} r \, dz \, dr \, d\theta = \frac{39}{4}\pi + \frac{1}{4}\pi\sqrt{21}$$

14. A ball of radius 12 is placed in a drill press and a hole of radius 4 is drilled out with the center of the hole a diameter of the ball. What is the volume of the material which remains?

Answer:

You know the formula for the volume of a sphere and so if you find out how much stuff is taken away, then it will be easy to find what is left. To find the volume of what is removed, it is easiest to use cylindrical coordinates. This volume is

$$\int_0^4 \int_0^{2\pi} \int_{-\sqrt{(144-r^2)}}^{\sqrt{(144-r^2)}} r \, dz \, d\theta \, dr = -\frac{4096}{3}\sqrt{2}\pi + 2304\pi.$$

Therefore, the volume of what remains is  $\frac{4}{3}\pi(12)^3$  minus the above. Thus the volume of what remains is

$$\frac{4096}{3}\sqrt{2}\pi.$$

15. A ball of radius 11 has density equal to  $\sqrt{x^2 + y^2 + z^2}$  in rectangular coordinates. The top of this ball is sliced off by a plane of the form  $z = 1$ . What is the mass of what remains?

Answer:

$$\begin{aligned} & \int_0^{2\pi} \int_0^{\arcsin(\frac{2}{11}\sqrt{30})} \int_0^{\sec \phi} \rho^3 \sin \phi \, d\rho \, d\phi \, d\theta + \int_0^{2\pi} \int_{\arcsin(\frac{2}{11}\sqrt{30})}^{\pi} \int_0^{11} \rho^3 \sin \phi \, d\rho \, d\phi \, d\theta \\ &= \frac{24623}{3}\pi \end{aligned}$$

16. Find  $\int \int_S \frac{y}{x} dV$  where  $S$  is described in polar coordinates as  $1 \leq r \leq 2$  and  $0 \leq \theta \leq \pi/4$ .

Answer:

Use  $x = r \cos \theta$  and  $y = r \sin \theta$ . Then the integral in polar coordinates is

$$\int_0^{\pi/4} \int_1^2 (r \tan \theta) dr d\theta = \frac{3}{4} \ln 2.$$

17. Find  $\int \int_S \left( \left( \frac{y}{x} \right)^2 + 1 \right) dV$  where  $S$  is given in polar coordinates as  $1 \leq r \leq 2$  and  $0 \leq \theta \leq \frac{1}{4}\pi$ .

Answer:

Use  $x = r \cos \theta$  and  $y = r \sin \theta$ . Then the integral in polar coordinates is

$$\int_0^{\frac{1}{4}\pi} \int_1^2 (1 + \tan^2 \theta) r dr d\theta.$$

18. Use polar coordinates to evaluate the following integral. Here  $S$  is given in terms of the polar coordinates.  $\int \int_S \sin(4x^2 + 4y^2) dV$  where  $r \leq 2$  and  $0 \leq \theta \leq \frac{1}{6}\pi$ .

Answer:

$$\int_0^{\frac{1}{6}\pi} \int_0^2 \sin(4r^2) r dr d\theta.$$

19. Find  $\int \int_S e^{2x^2+2y^2} dV$  where  $S$  is given in terms of the polar coordinates,  $r \leq 2$  and  $0 \leq \theta \leq \frac{1}{3}\pi$ .

Answer:

The integral is

$$\int_0^{\frac{1}{3}\pi} \int_0^2 r e^{2r^2} dr d\theta = \frac{1}{12} \pi (e^8 - 1).$$

20. Compute the volume of a sphere of radius  $R$  using cylindrical coordinates.

Answer:

Using cylindrical coordinates, the integral is  $\int_0^{2\pi} \int_0^R \int_{-\sqrt{R^2-r^2}}^{\sqrt{R^2-r^2}} r dz dr d\theta = \frac{4}{3} \pi R^3$ .

## 21.8 Exercises

- Find the area of the bounded region,  $R$ , determined by  $5x + y = 2$ ,  $5x + y = 8$ ,  $y = 2x$ , and  $y = 6x$ .
- Find the area of the bounded region,  $R$ , determined by  $y + 2x = 6$ ,  $y + 2x = 10$ ,  $y = 3x$ , and  $y = 4x$ .
- A solid,  $R$  is determined by  $3x + y = 2$ ,  $3x + y = 4$ ,  $y = 2x$ , and  $y = 6x$  and the density is  $\rho = x$ . Find the total mass of  $R$ .
- A solid,  $R$  is determined by  $4x + 2y = 5$ ,  $4x + 2y = 6$ ,  $y = 5x$ , and  $y = 7x$  and the density is  $\rho = y$ . Find the total mass of  $R$ .
- A solid,  $R$  is determined by  $3x + y = 3$ ,  $3x + y = 10$ ,  $y = 3x$ , and  $y = 5x$  and the density is  $\rho = y^{-1}$ . Find the total mass of  $R$ .

6. Find the volume of the region,  $E$ , bounded by the ellipsoid,  $\frac{1}{4}x^2 + y^2 + z^2 = 1$ .
7. Here are three vectors.  $(4, 1, 2)^T$ ,  $(5, 0, 2)^T$ , and  $(3, 1, 3)^T$ . These vectors determine a parallelepiped,  $R$ , which is occupied by a solid having density  $\rho = x$ . Find the mass of this solid.
8. Here are three vectors.  $(5, 1, 6)^T$ ,  $(6, 0, 6)^T$ , and  $(4, 1, 7)^T$ . These vectors determine a parallelepiped,  $R$ , which is occupied by a solid having density  $\rho = y$ . Find the mass of this solid.
9. Here are three vectors.  $(5, 2, 9)^T$ ,  $(6, 1, 9)^T$ , and  $(4, 2, 10)^T$ . These vectors determine a parallelepiped,  $R$ , which is occupied by a solid having density  $\rho = y + x$ . Find the mass of this solid.
10. Let  $D = \{(x, y) : x^2 + y^2 \leq 25\}$ . Find  $\int_D e^{25x^2 + 25y^2} dx dy$ .
11. Let  $D = \{(x, y) : x^2 + y^2 \leq 16\}$ . Find  $\int_D \cos(9x^2 + 9y^2) dx dy$ .
12. The ice cream in a sugar cone is described in spherical coordinates by  $\rho \in [0, 10]$ ,  $\phi \in [0, \frac{1}{3}\pi]$ ,  $\theta \in [0, 2\pi]$ . If the units are in centimeters, find the total volume in cubic centimeters of this ice cream.
13. Find the volume between  $z = 5 - x^2 - y^2$  and  $z = 2\sqrt{(x^2 + y^2)}$ .
14. A ball of radius 3 is placed in a drill press and a hole of radius 2 is drilled out with the center of the hole a diameter of the ball. What is the volume of the material which remains?
15. A ball of radius 9 has density equal to  $\sqrt{x^2 + y^2 + z^2}$  in rectangular coordinates. The top of this ball is sliced off by a plane of the form  $z = 2$ . What is the mass of what remains?
16. Find  $\int_S \frac{y}{x} dV$  where  $S$  is described in polar coordinates as  $1 \leq r \leq 2$  and  $0 \leq \theta \leq \pi/4$ .
17. Find  $\int_S \left( \left( \frac{y}{x} \right)^2 + 1 \right) dV$  where  $S$  is given in polar coordinates as  $1 \leq r \leq 2$  and  $0 \leq \theta \leq \frac{1}{6}\pi$ .
18. Use polar coordinates to evaluate the following integral. Here  $S$  is given in terms of the polar coordinates.  $\int_S \sin(2x^2 + 2y^2) dV$  where  $r \leq 2$  and  $0 \leq \theta \leq \frac{3}{2}\pi$ .
19. Find  $\int_S e^{2x^2 + 2y^2} dV$  where  $S$  is given in terms of the polar coordinates,  $r \leq 2$  and  $0 \leq \theta \leq \pi$ .
20. Compute the volume of a sphere of radius  $R$  using cylindrical coordinates.
21. In Example 21.6.6 on Page 541 check out all the details by working the integrals to be sure the steps are right.
22. What if the hollow sphere in Example 21.6.6 were in two dimensions and everything, including Newton's law still held? Would similar conclusions hold? Explain.
23. Fill in all details for the following argument that  $\int_0^\infty e^{-x^2} dx = \frac{1}{2}\sqrt{\pi}$ . Let  $I = \int_0^\infty e^{-x^2} dx$ . Then

$$I^2 = \int_0^\infty \int_0^\infty e^{-(x^2+y^2)} dx dy = \int_0^{\pi/2} \int_0^\infty r e^{-r^2} dr d\theta = \frac{1}{4}\pi$$

from which the result follows.

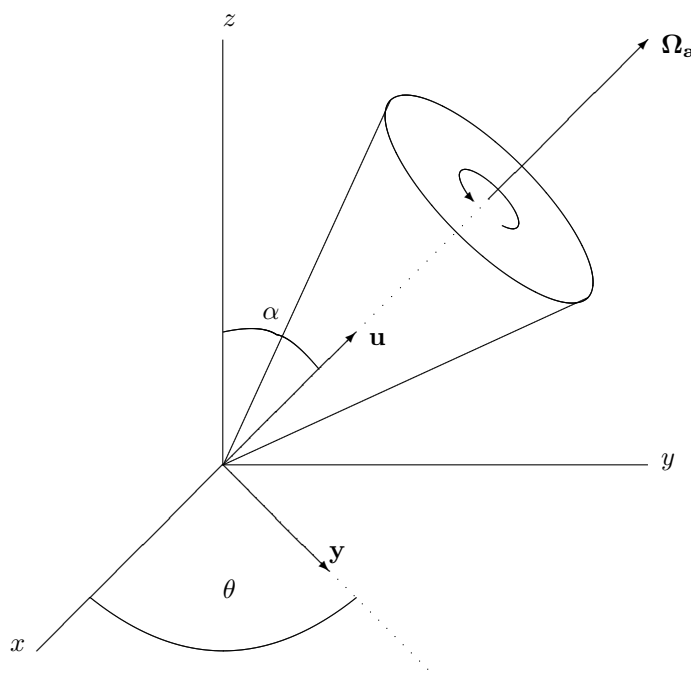
24. Show using Problem 23  $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$ , thus completing the discussion of the gamma function in the exercises on Page 256.
25. Let  $p, q > 0$  and define  $B(p, q) = \int_0^1 x^{p-1} (1-x)^{q-1}$ . Show  $\Gamma(p) \Gamma(q) = B(p, q) \Gamma(p+q)$ .  
**Hint:** It is fairly routine if you start with the left side and proceed to change variables.

## 21.9 The Moment Of Inertia

In order to appreciate the importance of this concept, it is necessary to discuss its physical significance.

### 21.9.1 The Spinning Top

To begin with consider a spinning top as illustrated in the following picture.



For the purpose of this discussion, consider the top as a large number of point masses,  $m_i$ , located at the positions,  $\mathbf{r}_i(t)$  for  $i = 1, 2, \dots, N$  and these masses are symmetrically arranged relative to the axis of the top. As the top spins, the axis of symmetry is observed to move around the  $z$  axis. This is called precession and you will see it occur whenever you spin a top. What is the speed of this precession? In other words, what is  $\theta'$ ?

Imagine a coordinate system which is fixed relative to the moving top. Thus in this coordinate system the points of the top are fixed. Let the standard unit vectors of the coordinate system moving with the top be denoted by  $\mathbf{i}(t), \mathbf{j}(t), \mathbf{k}(t)$  where  $\mathbf{i}(0), \mathbf{j}(0), \mathbf{k}(0)$  are the standard unit vectors corresponding to the coordinate axes indicated in the picture, also denoted by  $\mathbf{i}, \mathbf{j}, \mathbf{k}$ . Now take a vector,  $\mathbf{u}$  and suppose at time 0,

$$\mathbf{u} \equiv u^1 \mathbf{i}(0) + u^2 \mathbf{j}(0) + u^3 \mathbf{k}(0).$$

Let  $\mathbf{u}(t)$  be defined as the vector which has the same components with respect to  $\mathbf{i}(t), \mathbf{j}(t), \mathbf{k}(t)$  but at time  $t$ . Thus

$$\mathbf{u}(t) \equiv u^1 \mathbf{i}(t) + u^2 \mathbf{j}(t) + u^3 \mathbf{k}(t).$$

and the vector has changed although the components have not. Letting  $Q(t)\mathbf{u} \equiv \mathbf{u}(t)$ , it follows  $Q(t)$  is a linear transformation which preserves lengths. Therefore, from Lemma 19.5.3 on Page 429, There exists a vector  $\boldsymbol{\Omega}(t)$  such that

$$\mathbf{u}'(t) = \boldsymbol{\Omega}(t) \times \mathbf{u}(t).$$

The vector  $\boldsymbol{\Omega}_a$  shown in the picture is the vector for which

$$\mathbf{r}'_i(t) \equiv \boldsymbol{\Omega}_a \times \mathbf{r}_i(t)$$

is the velocity of the  $i^{th}$  point mass due to rotation about the axis of the top. Thus  $\boldsymbol{\Omega}(t) = \boldsymbol{\Omega}_a(t) + \boldsymbol{\Omega}_p(t)$  and it is assumed  $\boldsymbol{\Omega}_p(t)$  is very small relative to  $\boldsymbol{\Omega}_a$ . In other words, it is assumed the axis of the top moves very slowly relative to the speed of the points in the top which are spinning very fast around the axis of the top. The angular momentum,  $\mathbf{L}$  is defined by

$$\mathbf{L} \equiv \sum_{i=1}^N \mathbf{r}_i \times m_i \mathbf{v}_i \quad (21.14)$$

where  $\mathbf{v}_i$  equals the velocity of the  $i^{th}$  point mass. Thus  $\mathbf{v}_i = \boldsymbol{\Omega}(t) \times \mathbf{r}_i$  and from the above assumption,  $\mathbf{v}_i$  may be taken equal to  $\boldsymbol{\Omega}_a \times \mathbf{r}_i$ . Therefore,  $\mathbf{L}$  is essentially given by

$$\begin{aligned} \mathbf{L} &\equiv \sum_{i=1}^N m_i \mathbf{r}_i \times (\boldsymbol{\Omega}_a \times \mathbf{r}_i) \\ &= \sum_{i=1}^N m_i \left( |\mathbf{r}_i|^2 \boldsymbol{\Omega}_a - (\mathbf{r}_i \cdot \boldsymbol{\Omega}_a) \mathbf{r}_i \right). \end{aligned}$$

By symmetry of the top, this last expression equals a multiple of  $\boldsymbol{\Omega}_a$ . Thus  $\mathbf{L}$  is parallel to  $\boldsymbol{\Omega}_a$ . Also,

$$\begin{aligned} \mathbf{L} \cdot \boldsymbol{\Omega}_a &= \sum_{i=1}^N m_i \boldsymbol{\Omega}_a \cdot \mathbf{r}_i \times (\boldsymbol{\Omega}_a \times \mathbf{r}_i) \\ &= \sum_{i=1}^N m_i (\boldsymbol{\Omega}_a \times \mathbf{r}_i) \cdot (\boldsymbol{\Omega}_a \times \mathbf{r}_i) \\ &= \sum_{i=1}^N m_i |\boldsymbol{\Omega}_a \times \mathbf{r}_i|^2 = \sum_{i=1}^N m_i |\boldsymbol{\Omega}_a|^2 |\mathbf{r}_i|^2 \sin^2(\beta_i) \end{aligned}$$

where  $\beta_i$  denotes the angle between the position vector of the  $i^{th}$  point mass and the axis of the top. Since this expression is positive, this also shows  $\mathbf{L}$  has the same direction as  $\boldsymbol{\Omega}_a$ . Let  $\omega \equiv |\boldsymbol{\Omega}_a|$ . Then the above expression is of the form

$$\mathbf{L} \cdot \boldsymbol{\Omega}_a = I\omega^2,$$

where

$$I \equiv \sum_{i=1}^N m_i |\mathbf{r}_i|^2 \sin^2(\beta_i).$$



Thus, to get  $I$  you take the mass of the  $i^{th}$  point mass, multiply it by the square of its distance to the axis of the top and add all these up. This is defined as the moment of inertia of the top about the axis of the top. Letting  $\mathbf{u}$  denote a unit vector in the direction of the axis of the top, this implies

$$\mathbf{L} = I\omega\mathbf{u}. \quad (21.15)$$

Note the simple description of the angular momentum in terms of the moment of inertia. Referring to the above picture, define the vector,  $\mathbf{y}$  to be the projection of the vector,  $\mathbf{u}$  on the  $xy$  plane. Thus

$$\mathbf{y} = \mathbf{u} - (\mathbf{u} \cdot \mathbf{k})\mathbf{k}$$

and

$$(\mathbf{u} \cdot \mathbf{i}) = (\mathbf{y} \cdot \mathbf{i}) = \sin \alpha \cos \theta. \quad (21.16)$$

Now also from (21.14),

$$\begin{aligned} \frac{d\mathbf{L}}{dt} &= \sum_{i=1}^N m_i \overbrace{\mathbf{r}'_i \times \mathbf{v}_i}^{=0} + \mathbf{r}_i \times m_i \mathbf{v}'_i \\ &= \sum_{i=1}^N \mathbf{r}_i \times m_i \mathbf{v}'_i = - \sum_{i=1}^N \mathbf{r}_i \times m_i g \mathbf{k} \end{aligned}$$

where  $g$  is the acceleration of gravity. From (21.15), (21.16), and the above,

$$\begin{aligned} \frac{d\mathbf{L}}{dt} \cdot \mathbf{i} &= I\omega \left( \frac{d\mathbf{u}}{dt} \cdot \mathbf{i} \right) = I\omega \left( \frac{d\mathbf{y}}{dt} \cdot \mathbf{i} \right) \\ &= (-I\omega \sin \alpha \sin \theta) \theta' = - \sum_{i=1}^N \mathbf{r}_i \times m_i g \mathbf{k} \cdot \mathbf{i} \\ &= - \sum_{i=1}^N m_i g \mathbf{r}_i \cdot \mathbf{k} \times \mathbf{i} = - \sum_{i=1}^N m_i g \mathbf{r}_i \cdot \mathbf{j}. \end{aligned} \quad (21.17)$$

To simplify this further, recall the following definition of the center of mass.

**Definition 21.9.1** Define the total mass,  $M$  by

$$M = \sum_{i=1}^N m_i$$

and the center of mass,  $\mathbf{r}_0$  by

$$\mathbf{r}_0 \equiv \frac{\sum_{i=1}^N \mathbf{r}_i m_i}{M}. \quad (21.18)$$

In terms of the center of mass, the last expression equals

$$\begin{aligned} -Mg\mathbf{r}_0 \cdot \mathbf{j} &= -Mg(\mathbf{r}_0 - (\mathbf{r}_0 \cdot \mathbf{k})\mathbf{k} + (\mathbf{r}_0 \cdot \mathbf{k})\mathbf{k}) \cdot \mathbf{j} \\ &= -Mg(\mathbf{r}_0 - (\mathbf{r}_0 \cdot \mathbf{k})\mathbf{k}) \cdot \mathbf{j} \\ &= -Mg|\mathbf{r}_0 - (\mathbf{r}_0 \cdot \mathbf{k})\mathbf{k}| \cos \theta \\ &= -Mg|\mathbf{r}_0| \sin \alpha \cos \left( \frac{\pi}{2} - \theta \right). \end{aligned}$$

Note that by symmetry,  $\mathbf{r}_0(t)$  is on the axis of the top, is in the same direction as  $\mathbf{L}$ ,  $\mathbf{u}$ , and  $\boldsymbol{\Omega}_a$ , and also  $|\mathbf{r}_0|$  is independent of  $t$ . Therefore, from the second line of (21.17),

$$(-I\omega \sin \alpha \sin \theta) \theta' = -Mg |\mathbf{r}_0| \sin \alpha \sin \theta.$$

which shows

$$\theta' = \frac{Mg |\mathbf{r}_0|}{I\omega}. \quad (21.19)$$

From (21.19), the angular velocity of precession does not depend on  $\alpha$  in the picture. It also is slower when  $\omega$  is large and  $I$  is large.

The above discussion is a considerable simplification of the problem of a spinning top obtained from an assumption that  $\boldsymbol{\Omega}_a$  is approximately equal to  $\boldsymbol{\Omega}$ . It also leaves out all considerations of friction and the observation that the axis of symmetry wobbles. This wobbling is called nutation. The full mathematical treatment of this problem involves the Euler angles and some fairly complicated differential equations obtained using techniques discussed in advanced physics classes. Lagrange studied these types of problems back in the 1700's.

### 21.9.2 Kinetic Energy

The next problem is that of understanding the total kinetic energy of a collection of moving point masses. Consider a possibly large number of point masses,  $m_i$  located at the positions  $\mathbf{r}_i$  for  $i = 1, 2, \dots, N$ . Thus the velocity of the  $i^{\text{th}}$  point mass is  $\mathbf{r}'_i = \mathbf{v}_i$ . The kinetic energy of the mass  $m_i$  is defined by

$$\frac{1}{2} m_i |\mathbf{r}'_i|^2.$$

(This is a very good time to review the presentation on kinetic energy given on Page 363.) The total kinetic energy of the collection of masses is then

$$E = \sum_{i=1}^N \frac{1}{2} m_i |\mathbf{r}'_i|^2. \quad (21.20)$$

As these masses move about, so does the center of mass,  $\mathbf{r}_0$ . Thus  $\mathbf{r}_0$  is a function of  $t$  just as the other  $\mathbf{r}_i$ . From (21.20) the total kinetic energy is

$$\begin{aligned} E &= \sum_{i=1}^N \frac{1}{2} m_i |\mathbf{r}'_i - \mathbf{r}'_0 + \mathbf{r}'_0|^2 \\ &= \sum_{i=1}^N \frac{1}{2} m_i \left[ |\mathbf{r}'_i - \mathbf{r}'_0|^2 + |\mathbf{r}'_0|^2 + 2(\mathbf{r}'_i - \mathbf{r}'_0) \cdot \mathbf{r}'_0 \right]. \end{aligned} \quad (21.21)$$

Now

$$\begin{aligned} \sum_{i=1}^N m_i (\mathbf{r}'_i - \mathbf{r}'_0) \cdot \mathbf{r}'_0 &= \left( \sum_{i=1}^N m_i (\mathbf{r}_i - \mathbf{r}_0) \right)' \cdot \mathbf{r}'_0 \\ &= 0 \end{aligned}$$

because from (21.18)

$$\begin{aligned} \sum_{i=1}^N m_i (\mathbf{r}_i - \mathbf{r}_0) &= \sum_{i=1}^N m_i \mathbf{r}_i - \sum_{i=1}^N m_i \mathbf{r}_0 \\ &= \sum_{i=1}^N m_i \mathbf{r}_i - \sum_{i=1}^N m_i \left( \frac{\sum_{i=1}^N \mathbf{r}_i m_i}{\sum_{i=1}^N m_i} \right) = \mathbf{0}. \end{aligned}$$

Let  $M \equiv \sum_{i=1}^N m_i$  be the total mass. Then (21.21) reduces to

$$\begin{aligned} E &= \sum_{i=1}^N \frac{1}{2} m_i \left[ |\mathbf{r}'_i - \mathbf{r}'_0|^2 + |\mathbf{r}'_0|^2 \right] \\ &= \frac{1}{2} M |\mathbf{r}'_0|^2 + \sum_{i=1}^N \frac{1}{2} m_i |\mathbf{r}'_i - \mathbf{r}'_0|^2. \end{aligned} \quad (21.22)$$

The first term is just the kinetic energy of a point mass equal to the sum of all the masses involved, located at the center of mass of the system of masses while the second term represents kinetic energy which comes from the relative velocities of the masses taken with respect to the center of mass. It is this term which is considered more carefully in the case of rigid body motion.

**Definition 21.9.2** *The collection of masses,  $\{m_i\}$  located at the positions,  $\mathbf{r}_i$  is said to undergo rigid body motion if whenever  $i \neq j$ ,  $|\mathbf{r}_i(t) - \mathbf{r}_j(t)|$  is a constant which does not depend on  $t$ . Thus the distance between any pair of these point masses does not change.*

A fun experiment is to take a hard boiled egg and spin it and then take a raw egg and give it a spin. You will certainly feel a big difference in the way the two eggs respond. Incidentally, this is a good way to tell whether the egg has been hard boiled or is raw and can be used to prevent messiness which could occur if you think it is hard boiled and it really isn't.

Now let  $\mathbf{e}_1(t)$ ,  $\mathbf{e}_2(t)$ , and  $\mathbf{e}_3(t)$  be an orthonormal set of vectors which is fixed in the body undergoing rigid body motion. This means that  $\mathbf{r}_i(t) - \mathbf{r}_0(t)$  has components which are constant in  $t$  with respect to the vectors,  $\mathbf{e}_i(t)$ . Now let  $Q(t)$  be a linear transformation defined by

$$Q(t) \mathbf{u} \equiv u_1 \mathbf{e}_1(t) + u_2 \mathbf{e}_2(t) + u_3 \mathbf{e}_3(t)$$

where

$$\mathbf{u} \equiv u_1 \mathbf{e}_1(0) + u_2 \mathbf{e}_2(0) + u_3 \mathbf{e}_3(0).$$

Thus

$$\mathbf{r}_i(t) - \mathbf{r}_0(t) = Q(t) (\mathbf{r}_i(0) - \mathbf{r}_0(0)).$$

Then as in the discussion of the Coriolis force presented earlier, Lemma 19.5.3 on Page 429 implies  $Q'(t) Q(t)^T$  is of the form

$$\begin{pmatrix} 0 & -\omega_3(t) & \omega_2(t) \\ \omega_3(t) & 0 & -\omega_1(t) \\ -\omega_2(t) & \omega_1(t) & 0 \end{pmatrix}$$

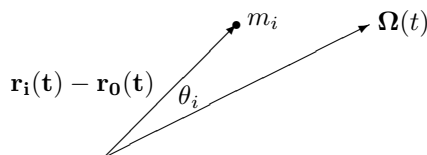
and there exists a vector,  $\boldsymbol{\Omega}(t)$  which does not depend on  $i$  such that

$$\mathbf{r}'_i(t) - \mathbf{r}'_0(t) = \boldsymbol{\Omega}(t) \times (\mathbf{r}_i(t) - \mathbf{r}_0(t)).$$

Now using this in (21.22),

$$\begin{aligned} E &= \frac{1}{2} M |\mathbf{r}'_0|^2 + \sum_{i=1}^N \frac{1}{2} m_i |\boldsymbol{\Omega}(t) \times (\mathbf{r}_i(t) - \mathbf{r}_0(t))|^2 \\ &= \frac{1}{2} M |\mathbf{r}'_0|^2 + \frac{1}{2} \left( \sum_{i=1}^N m_i |\mathbf{r}_i(t) - \mathbf{r}_0(t)|^2 \sin^2 \theta_i \right) |\boldsymbol{\Omega}(t)|^2 \\ &= \frac{1}{2} M |\mathbf{r}'_0|^2 + \frac{1}{2} \left( \sum_{i=1}^N m_i |\mathbf{r}_i(0) - \mathbf{r}_0(0)|^2 \sin^2 \theta_i \right) |\boldsymbol{\Omega}(t)|^2 \end{aligned}$$

where  $\theta_i$  is the angle between  $\boldsymbol{\Omega}(t)$  and the vector,  $\mathbf{r}_i(t) - \mathbf{r}_0(t)$ . Therefore,  $|\mathbf{r}_i(t) - \mathbf{r}_0(t)| \sin \theta_i$  is the distance between the point mass,  $m_i$  located at  $\mathbf{r}_i$  and a line through the center of mass,  $\mathbf{r}_0$  with direction,  $\boldsymbol{\Omega}$  as indicated in the following picture.



Thus the expression,  $\sum_{i=1}^N m_i |\mathbf{r}_i(0) - \mathbf{r}_0(0)|^2 \sin^2 \theta_i$  plays the role of a mass in the definition of kinetic energy except instead of the speed, substitute the angular speed,  $|\boldsymbol{\Omega}(t)|$ . It is this expression which is called the moment of inertia about the line whose direction is  $\boldsymbol{\Omega}(t)$ .

In both of these examples, the center of mass and the moment of inertia occurred in a natural way.

### 21.9.3 Finding The Moment Of Inertia And Center Of Mass

The methods used to evaluate multiple integrals make possible the determination of centers of mass and moments of inertia. In the case of a solid material rather than finitely many point masses, you replace the sums with integrals. The sums are essentially approximations of the integrals which result.

**Example 21.9.3** *Let a solid occupy the three dimensional region  $R$  and suppose the density is  $\rho$ . What is the moment of inertia of this solid about the  $z$  axis? What is the center of mass?*

Here the little masses would be of the form  $\rho(\mathbf{x}) dV$  where  $\mathbf{x}$  is a point of  $R$ . Therefore, the contribution of this mass to the moment of inertia would be

$$(x^2 + y^2) \rho(\mathbf{x}) dV$$

where the Cartesian coordinates of the point  $\mathbf{x}$  are  $(x, y, z)$ . Then summing these up as an integral, yields the following for the moment of inertia.

$$\int_R (x^2 + y^2) \rho(\mathbf{x}) dV. \quad (21.23)$$

To find the center of mass, sum up  $\mathbf{r} \rho dV$  for the points in  $R$  and divide by the total mass. In Cartesian coordinates, where  $\mathbf{r} = (x, y, z)$ , this means to sum up vectors of the form  $(x \rho dV, y \rho dV, z \rho dV)$  and divide by the total mass. Thus the Cartesian coordinates of the center of mass are

$$\left( \frac{\int_R x \rho dV}{\int_R \rho dV}, \frac{\int_R y \rho dV}{\int_R \rho dV}, \frac{\int_R z \rho dV}{\int_R \rho dV} \right) \equiv \frac{\int_R \mathbf{r} \rho dV}{\int_R \rho dV}.$$

Here is a specific example.

**Example 21.9.4** *Find the moment of inertia about the  $z$  axis and center of mass of the solid which occupies the region,  $R$  defined by  $9 - (x^2 + y^2) \geq z \geq 0$  if the density is  $\rho(x, y, z) = \sqrt{x^2 + y^2}$ .*

This moment of inertia is  $\int_R (x^2 + y^2) \sqrt{x^2 + y^2} dV$  and the easiest way to find this integral is to use cylindrical coordinates. Thus the answer is

$$\int_0^{2\pi} \int_0^3 \int_0^{9-r^2} r^3 dz dr d\theta = \frac{8748}{35} \pi.$$

To find the center of mass, note the  $x$  and  $y$  coordinates of the center of mass,

$$\frac{\int_R x \rho dV}{\int_R \rho dV}, \frac{\int_R y \rho dV}{\int_R \rho dV}$$

both equal zero because the above shape is symmetric about the  $z$  axis and  $\rho$  is also symmetric in its values. Thus  $x\rho dV$  will cancel with  $-x\rho dV$  and a similar conclusion will hold for the  $y$  coordinate. It only remains to find the  $z$  coordinate of the center of mass,  $\bar{z}$ . In polar coordinates,  $\rho = r$  and so,

$$\bar{z} = \frac{\int_R z \rho dV}{\int_R \rho dV} = \frac{\int_0^{2\pi} \int_0^3 \int_0^{9-r^2} z r^2 dz dr d\theta}{\int_0^{2\pi} \int_0^3 \int_0^{9-r^2} r^2 dz dr d\theta} = \frac{18}{7}.$$

Thus the center of mass will be  $(0, 0, \frac{18}{7})$ .

## 21.10 Exercises

1. Let  $R$  denote the finite region bounded by  $z = 4 - x^2 - y^2$  and the  $xy$  plane. Find  $z_c$ , the  $z$  coordinate of the center of mass if the density,  $\sigma$  is a constant.
2. Let  $B$  be a solid ball of constant density and radius  $R$ . Find the moment of inertia about a line through a diameter of the sphere. You should get  $\frac{2}{5}R^2M$ .
3. Let  $C$  be a solid cylinder of constant density and radius  $R$ . Find the moment of inertia about the axis of the cylinder  
You should get  $\frac{1}{2}R^2M$ .
4. Let  $C$  be a solid cylinder of constant density and radius  $R$  and mass  $M$  and let  $B$  be a solid ball of radius  $R$  and mass  $M$ . The cylinder and the sphere are placed on the top of an inclined plane and allowed to roll to the bottom. Which one will arrive first and why?
5. Suppose a solid of mass  $M$  occupying the region,  $B$  has moment of inertia,  $I_l$  about a line,  $l$  which passes through the center of mass of  $M$  and let  $l_1$  be another line parallel to  $l$  and at a distance of  $a$  from  $l$ . Then the parallel axis theorem states  $I_{l_1} = I_l + a^2M$ . Prove the parallel axis theorem. **Hint:** Choose axes such that the  $z$  axis is  $l$  and  $l_1$  passes through the point  $(a, 0)$  in the  $xy$  plane.
6. Using the parallel axis theorem find the moment of inertia of a solid ball of radius  $R$  and mass  $M$  about an axis located at a distance of  $a$  from the center of the ball. Your answer should be  $Ma^2 + \frac{2}{5}MR^2$ .
7. Find the moment of inertia of a solid thin rod of length  $l$ , mass  $M$ , and constant density about an axis through the center of the rod perpendicular to the axis of the rod. You should get  $\frac{1}{12}l^2M$ .

8. Using the parallel axis theorem, find the moment of inertia of a solid thin rod of length  $l$ , mass  $M$ , and constant density about an axis through an end of the rod perpendicular to the axis of the rod. You should get  $\frac{1}{3}l^2M$ .
9. Let the angle between the  $z$  axis and the sides of a right circular cone be  $\alpha$ . Also assume the height of this cone is  $h$ . Find the  $z$  coordinate of the center of mass of this cone assuming the density is constant. **Hint:** Use spherical coordinates. You should get  $(3/4)h$ .
10. Let  $R$  denote the part of the solid ball,  $x^2 + y^2 + z^2 \leq R^2$  which lies in the first octant. That is  $x, y, z \geq 0$ . Find the coordinates of the center of mass if the density is constant. Your answer for one of the coordinates for the center of mass should be  $(3/8)R$ .
11. Show that in general for  $\mathbf{L}$  angular momentum,

$$\frac{d\mathbf{L}}{dt} = \mathbf{\Gamma}$$

where  $\mathbf{\Gamma}$  is the total torque,

$$\mathbf{\Gamma} \equiv \sum \mathbf{r}_i \times \mathbf{F}_i$$

where  $\mathbf{F}_i$  is the force on the  $i^{th}$  point mass.

## 21.11 Theory Of The Riemann Integral

The definition of the Riemann integral of a function of  $n$  variables uses the following definition.

**Definition 21.11.1** For  $i = 1, \dots, n$ , let  $\{\alpha_k^i\}_{k=-\infty}^{\infty}$  be points on  $\mathbb{R}$  which satisfy

$$\lim_{k \rightarrow \infty} \alpha_k^i = \infty, \quad \lim_{k \rightarrow -\infty} \alpha_k^i = -\infty, \quad \alpha_k^i < \alpha_{k+1}^i. \quad (21.24)$$

For such sequences, define a grid on  $\mathbb{R}^n$  denoted by  $\mathcal{G}$  or  $\mathcal{F}$  as the collection of boxes of the form

$$Q = \prod_{i=1}^n [\alpha_{j_i}^i, \alpha_{j_i+1}^i]. \quad (21.25)$$

If  $\mathcal{G}$  is a grid,  $\mathcal{F}$  is called a refinement of  $\mathcal{G}$  if every box of  $\mathcal{G}$  is the union of boxes of  $\mathcal{F}$ .

**Lemma 21.11.2** If  $\mathcal{G}$  and  $\mathcal{F}$  are two grids, they have a common refinement, denoted here by  $\mathcal{G} \vee \mathcal{F}$ .

**Proof:** Let  $\{\alpha_k^i\}_{k=-\infty}^{\infty}$  be the sequences used to construct  $\mathcal{G}$  and let  $\{\beta_k^i\}_{k=-\infty}^{\infty}$  be the sequence used to construct  $\mathcal{F}$ . Now let  $\{\gamma_k^i\}_{k=-\infty}^{\infty}$  denote the union of  $\{\alpha_k^i\}_{k=-\infty}^{\infty}$  and  $\{\beta_k^i\}_{k=-\infty}^{\infty}$ . It is necessary to show that for each  $i$  these points can be arranged in order. To do so, let  $\gamma_0^i \equiv \alpha_0^i$ . Now if

$$\gamma_{-j}^i, \dots, \gamma_0^i, \dots, \gamma_j^i$$

have been chosen such that they are in order and all distinct, let  $\gamma_{j+1}^i$  be the first element of

$$\{\alpha_k^i\}_{k=-\infty}^{\infty} \cup \{\beta_k^i\}_{k=-\infty}^{\infty} \quad (21.26)$$

which is larger than  $\gamma_j^i$  and let  $\gamma_{-(j+1)}^i$  be the last element of (21.26) which is strictly smaller than  $\gamma_{-j}^i$ . The assumption (21.24) insures such a first and last element exists. Now let the grid  $\mathcal{G} \vee \mathcal{F}$  consist of boxes of the form

$$Q \equiv \prod_{i=1}^n [\gamma_{j_i}^i, \gamma_{j_i+1}^i].$$

The Riemann integral is only defined for functions,  $f$  which are bounded and are equal to zero out of some bounded set,  $D$ . In what follows  $f$  will always be such a function.

**Definition 21.11.3** Let  $f$  be a bounded function which equals zero off a bounded set,  $D$ , and let  $\mathcal{G}$  be a grid. For  $Q \in \mathcal{G}$ , define

$$M_Q(f) \equiv \sup \{f(\mathbf{x}) : \mathbf{x} \in Q\}, \quad m_Q(f) \equiv \inf \{f(\mathbf{x}) : \mathbf{x} \in Q\}. \quad (21.27)$$

Also define for  $Q$  a box, the volume of  $Q$ , denoted by  $v(Q)$  by

$$v(Q) \equiv \prod_{i=1}^n (b_i - a_i), \quad Q \equiv \prod_{i=1}^n [a_i, b_i].$$

Now define upper sums,  $\mathcal{U}_{\mathcal{G}}(f)$  and lower sums,  $\mathcal{L}_{\mathcal{G}}(f)$  with respect to the indicated grid, by the formulas

$$\mathcal{U}_{\mathcal{G}}(f) \equiv \sum_{Q \in \mathcal{G}} M_Q(f) v(Q), \quad \mathcal{L}_{\mathcal{G}}(f) \equiv \sum_{Q \in \mathcal{G}} m_Q(f) v(Q).$$

A function of  $n$  variables is Riemann integrable when there is a unique number between all the upper and lower sums. This number is the value of the integral.

Note that in this definition,  $M_Q(f) = m_Q(f) = 0$  for all but finitely many  $Q \in \mathcal{G}$  so there are no convergence questions to be considered here.

**Lemma 21.11.4** If  $\mathcal{F}$  is a refinement of  $\mathcal{G}$  then

$$\mathcal{U}_{\mathcal{G}}(f) \geq \mathcal{U}_{\mathcal{F}}(f), \quad \mathcal{L}_{\mathcal{G}}(f) \leq \mathcal{L}_{\mathcal{F}}(f).$$

Also if  $\mathcal{F}$  and  $\mathcal{G}$  are two grids,

$$\mathcal{L}_{\mathcal{G}}(f) \leq \mathcal{U}_{\mathcal{F}}(f).$$

**Proof:** For  $P \in \mathcal{G}$  let  $\hat{P}$  denote the set,

$$\{Q \in \mathcal{F} : Q \subseteq P\}.$$

Then  $P = \cup \hat{P}$  and

$$\begin{aligned} \mathcal{L}_{\mathcal{F}}(f) &\equiv \sum_{Q \in \mathcal{F}} m_Q(f) v(Q) = \sum_{P \in \mathcal{G}} \sum_{Q \in \hat{P}} m_Q(f) v(Q) \\ &\geq \sum_{P \in \mathcal{G}} m_P(f) \sum_{Q \in \hat{P}} v(Q) = \sum_{P \in \mathcal{G}} m_P(f) v(P) \equiv \mathcal{L}_{\mathcal{G}}(f). \end{aligned}$$

Similarly, the other inequality for the upper sums is valid.

To verify the last assertion of the lemma, use Lemma 21.11.2 to write

$$\mathcal{L}_{\mathcal{G}}(f) \leq \mathcal{L}_{\mathcal{G} \vee \mathcal{F}}(f) \leq \mathcal{U}_{\mathcal{G} \vee \mathcal{F}}(f) \leq \mathcal{U}_{\mathcal{F}}(f).$$

This proves the lemma.

This lemma makes it possible to define the Riemann integral.

**Definition 21.11.5** Define an upper and a lower integral as follows.

$$\bar{I}(f) \equiv \inf \{ \mathcal{U}_{\mathcal{G}}(f) : \mathcal{G} \text{ is a grid} \},$$

$$\underline{I}(f) \equiv \sup \{ \mathcal{L}_{\mathcal{G}}(f) : \mathcal{G} \text{ is a grid} \}.$$

**Lemma 21.11.6**  $\bar{I}(f) \geq \underline{I}(f)$ .

**Proof:** From Lemma 21.11.4 it follows for any two grids  $\mathcal{G}$  and  $\mathcal{F}$ ,

$$\mathcal{L}_{\mathcal{G}}(f) \leq \mathcal{U}_{\mathcal{F}}(f).$$

Therefore, taking the supremum for all grids on the left in this inequality,

$$\underline{I}(f) \leq \mathcal{U}_{\mathcal{F}}(f)$$

for all grids  $\mathcal{F}$ . Taking the infimum in this inequality, yields the conclusion of the lemma.

**Definition 21.11.7** A bounded function,  $f$  which equals zero off a bounded set,  $D$ , is said to be Riemann integrable, written as  $f \in \mathcal{R}(\mathbb{R}^n)$  exactly when  $\underline{I}(f) = \bar{I}(f)$ . In this case define

$$\int f \, dV \equiv \int f \, dx = \bar{I}(f) = \underline{I}(f).$$

As in the case of integration of functions of one variable, one obtains the Riemann criterion which is stated as the following theorem.

**Theorem 21.11.8** (Riemann criterion)  $f \in \mathcal{R}(\mathbb{R}^n)$  if and only if for all  $\varepsilon > 0$  there exists a grid  $\mathcal{G}$  such that

$$\mathcal{U}_{\mathcal{G}}(f) - \mathcal{L}_{\mathcal{G}}(f) < \varepsilon.$$

**Proof:** If  $f \in \mathcal{R}(\mathbb{R}^n)$ , then  $\bar{I}(f) = \underline{I}(f)$  and so there exist grids  $\mathcal{G}$  and  $\mathcal{F}$  such that

$$\mathcal{U}_{\mathcal{G}}(f) - \mathcal{L}_{\mathcal{F}}(f) \leq \bar{I}(f) + \frac{\varepsilon}{2} - \left( \underline{I}(f) - \frac{\varepsilon}{2} \right) = \varepsilon.$$

Then letting  $\mathcal{H} = \mathcal{G} \vee \mathcal{F}$ , Lemma 21.11.4 implies

$$\mathcal{U}_{\mathcal{H}}(f) - \mathcal{L}_{\mathcal{H}}(f) \leq \mathcal{U}_{\mathcal{G}}(f) - \mathcal{L}_{\mathcal{F}}(f) < \varepsilon.$$

Conversely, if for all  $\varepsilon > 0$  there exists  $\mathcal{G}$  such that

$$\mathcal{U}_{\mathcal{G}}(f) - \mathcal{L}_{\mathcal{G}}(f) < \varepsilon,$$

then

$$\bar{I}(f) - \underline{I}(f) \leq \mathcal{U}_{\mathcal{G}}(f) - \mathcal{L}_{\mathcal{G}}(f) < \varepsilon.$$

Since  $\varepsilon > 0$  is arbitrary, this proves the theorem.

### 21.11.1 Basic Properties

It is important to know that certain combinations of Riemann integrable functions are Riemann integrable. The following theorem will include all the important cases.

**Theorem 21.11.9** Let  $f, g \in \mathcal{R}(\mathbb{R}^n)$  and let  $\phi : K \rightarrow \mathbb{R}$  be continuous where  $K$  is a compact set in  $\mathbb{R}^2$  containing  $f(\mathbb{R}^n) \times g(\mathbb{R}^n)$ . Also suppose that  $\phi(0, 0) = 0$ . Then defining

$$h(\mathbf{x}) \equiv \phi(f(\mathbf{x}), g(\mathbf{x})),$$

it follows that  $h$  is also in  $\mathcal{R}(\mathbb{R}^n)$ .



**Proof:** Let  $\varepsilon > 0$  and let  $\delta_1 > 0$  be such that if  $(y_i, z_i), i = 1, 2$  are points in  $K$ , such that  $|z_1 - z_2| \leq \delta_1$  and  $|y_1 - y_2| \leq \delta_1$ , then

$$|\phi(y_1, z_1) - \phi(y_2, z_2)| < \varepsilon.$$

Let  $0 < \delta < \min(\delta_1, \varepsilon, 1)$ . Let  $\mathcal{G}$  be a grid with the property that for  $Q \in \mathcal{G}$ , the diameter of  $Q$  is less than  $\delta$  and also for  $k = f, g$ ,

$$\mathcal{U}_{\mathcal{G}}(k) - \mathcal{L}_{\mathcal{G}}(k) < \delta^2. \quad (21.28)$$

Then defining for  $k = f, g$ ,

$$\mathcal{P}_k \equiv \{Q \in \mathcal{G} : M_Q(k) - m_Q(k) > \delta\},$$

it follows

$$\begin{aligned} \delta^2 &> \sum_{Q \in \mathcal{G}} (M_Q(k) - m_Q(k)) v(Q) \geq \\ &\sum_{\mathcal{P}_k} (M_Q(k) - m_Q(k)) v(Q) \geq \delta \sum_{\mathcal{P}_k} v(Q) \end{aligned}$$

and so for  $k = f, g$ ,

$$\varepsilon > \delta > \sum_{\mathcal{P}_k} v(Q). \quad (21.29)$$

Suppose for  $k = f, g$ ,

$$M_Q(k) - m_Q(k) \leq \delta.$$

Then if  $\mathbf{x}_1, \mathbf{x}_2 \in Q$ ,

$$|f(\mathbf{x}_1) - f(\mathbf{x}_2)| < \delta, \text{ and } |g(\mathbf{x}_1) - g(\mathbf{x}_2)| < \delta.$$

Therefore,

$$|h(\mathbf{x}_1) - h(\mathbf{x}_2)| \equiv |\phi(f(\mathbf{x}_1), g(\mathbf{x}_1)) - \phi(f(\mathbf{x}_2), g(\mathbf{x}_2))| < \varepsilon$$

and it follows that

$$|M_Q(h) - m_Q(h)| \leq \varepsilon.$$

Now let

$$\mathcal{S} \equiv \{Q \in \mathcal{G} : 0 < M_Q(k) - m_Q(k) \leq \delta, k = f, g\}.$$

Thus the union of the boxes in  $\mathcal{S}$  is contained in some large box,  $R$ , which depends only on  $f$  and  $g$  and also, from the assumption that  $\phi(0, 0) = 0$ ,  $M_Q(h) - m_Q(h) = 0$  unless  $Q \subseteq R$ . Then

$$\begin{aligned} \mathcal{U}_{\mathcal{G}}(h) - \mathcal{L}_{\mathcal{G}}(h) &\leq \sum_{Q \in \mathcal{P}_f} (M_Q(h) - m_Q(h)) v(Q) + \\ &\sum_{Q \in \mathcal{P}_g} (M_Q(h) - m_Q(h)) v(Q) + \sum_{Q \in \mathcal{S}} \delta v(Q). \end{aligned}$$

Now since  $K$  is compact, it follows  $\phi(K)$  is bounded and so there exists a constant,  $C$ , depending only on  $h$  and  $\phi$  such that  $M_Q(h) - m_Q(h) < C$ . Therefore, the above inequality implies

$$\mathcal{U}_{\mathcal{G}}(h) - \mathcal{L}_{\mathcal{G}}(h) \leq C \sum_{Q \in \mathcal{P}_f} v(Q) + C \sum_{Q \in \mathcal{P}_g} v(Q) + \sum_{Q \in \mathcal{S}} \delta v(Q),$$

which by (21.29) implies

$$\mathcal{U}_{\mathcal{G}}(h) - \mathcal{L}_{\mathcal{G}}(h) \leq 2C\varepsilon + \delta v(R) \leq 2C\varepsilon + \varepsilon v(R).$$

Since  $\varepsilon$  is arbitrary, the Riemann criterion is satisfied and so  $h \in \mathcal{R}(\mathbb{R}^n)$ .

**Corollary 21.11.10** *Let  $f, g \in \mathcal{R}(\mathbb{R}^n)$  and let  $a, b \in \mathbb{R}$ . Then  $af + bg$ ,  $fg$ , and  $|f|$  are all in  $\mathcal{R}(\mathbb{R}^n)$ . Also,*

$$\int_{\mathbb{R}^n} (af + bg) dx = a \int_{\mathbb{R}^n} f dx + b \int_{\mathbb{R}^n} g dx, \quad (21.30)$$

and

$$\int |f| dx \geq \left| \int f dx \right|. \quad (21.31)$$

**Proof:** Each of the combinations of functions described above is Riemann integrable by Theorem 21.11.9. For example, to see  $af + bg \in \mathcal{R}(\mathbb{R}^n)$  consider  $\phi(y, z) \equiv ay + bz$ . This is clearly a continuous function of  $(y, z)$  such that  $\phi(0, 0) = 0$ . To obtain  $|f| \in \mathcal{R}(\mathbb{R}^n)$ , let  $\phi(y, z) \equiv |y|$ . It remains to verify the formulas. To do so, let  $\mathcal{G}$  be a grid with the property that for  $k = f, g, |f|$  and  $af + bg$ ,

$$\mathcal{U}_{\mathcal{G}}(k) - \mathcal{L}_{\mathcal{G}}(k) < \varepsilon. \quad (21.32)$$

Consider (21.30). For each  $Q \in \mathcal{G}$  pick a point in  $Q$ ,  $\mathbf{x}_Q$ . Then

$$\sum_{Q \in \mathcal{G}} k(\mathbf{x}_Q) v(Q) \in [\mathcal{L}_{\mathcal{G}}(k), \mathcal{U}_{\mathcal{G}}(k)]$$

and so

$$\left| \int k dx - \sum_{Q \in \mathcal{G}} k(\mathbf{x}_Q) v(Q) \right| < \varepsilon.$$

Consequently, since

$$\begin{aligned} & \sum_{Q \in \mathcal{G}} (af + bg)(\mathbf{x}_Q) v(Q) \\ &= a \sum_{Q \in \mathcal{G}} f(\mathbf{x}_Q) v(Q) + b \sum_{Q \in \mathcal{G}} g(\mathbf{x}_Q) v(Q), \end{aligned}$$

it follows

$$\begin{aligned} & \left| \int (af + bg) dx - a \int f dx - b \int g dx \right| \leq \\ & \left| \int (af + bg) dx - \sum_{Q \in \mathcal{G}} (af + bg)(\mathbf{x}_Q) v(Q) \right| + \\ & \left| a \sum_{Q \in \mathcal{G}} f(\mathbf{x}_Q) v(Q) - a \int f dx \right| + \left| b \sum_{Q \in \mathcal{G}} g(\mathbf{x}_Q) v(Q) - b \int g dx \right| \\ & \leq \varepsilon + |a| \varepsilon + |b| \varepsilon. \end{aligned}$$

Since  $\varepsilon$  is arbitrary, this establishes Formula (21.30) and shows the integral is linear.

It remains to establish the inequality (21.31). By (21.32), and the triangle inequality for sums,

$$\begin{aligned} \int |f| dx + \varepsilon & \geq \sum_{Q \in \mathcal{G}} |f(\mathbf{x}_Q)| v(Q) \geq \\ & \geq \left| \sum_{Q \in \mathcal{G}} f(\mathbf{x}_Q) v(Q) \right| \geq \left| \int f dx \right| - \varepsilon. \end{aligned}$$

Then since  $\varepsilon$  is arbitrary, this establishes the desired inequality. This proves the corollary.

Which functions are in  $\mathcal{R}(\mathbb{R}^n)$ ? Begin with step functions defined below.

**Definition 21.11.11** If

$$Q \equiv \prod_{i=1}^n [a_i, b_i]$$

is a box, define  $\text{int}(Q)$  as

$$\text{int}(Q) \equiv \prod_{i=1}^n (a_i, b_i).$$

$f$  is called a step function if there is a grid,  $\mathcal{G}$  such that  $f$  is constant on  $\text{int}(Q)$  for each  $Q \in \mathcal{G}$ ,  $f$  is bounded, and  $f(\mathbf{x}) = 0$  for all  $\mathbf{x}$  outside some bounded set.

The next corollary states that step functions are in  $\mathcal{R}(\mathbb{R}^n)$  and shows the expected formula for the integral is valid.

**Corollary 21.11.12** Let  $\mathcal{G}$  be a grid and let  $f$  be a step function such that  $f = f_Q$  on  $\text{int}(Q)$  for each  $Q \in \mathcal{G}$ . Then  $f \in \mathcal{R}(\mathbb{R}^n)$  and

$$\int f dx = \sum_{Q \in \mathcal{G}} f_Q v(Q).$$

**Proof:** Let  $Q$  be a box of  $\mathcal{G}$ ,

$$Q \equiv \prod_{i=1}^n [\alpha_{j_i}^i, \alpha_{j_i+1}^i],$$

and suppose  $g$  is a bounded function,  $|g(\mathbf{x})| \leq C$ , and  $g = 0$  off  $Q$ , and  $g = 1$  on  $\text{int}(Q)$ . Thus,  $g$  is the simplest sort of step function. Refine  $\mathcal{G}$  by including the extra points,

$$\alpha_{j_i}^i + \eta \text{ and } \alpha_{j_i+1}^i - \eta$$

for each  $i = 1, \dots, n$ . Here  $\eta$  is small enough that for each  $i$ ,  $\alpha_{j_i}^i + \eta < \alpha_{j_i+1}^i - \eta$ . Also let  $L$  denote the largest of the lengths of the sides of  $Q$ . Let  $\mathcal{F}$  be this refined grid and denote by  $Q_\eta$  the box

$$\prod_{i=1}^n [\alpha_{j_i}^i + \eta, \alpha_{j_i+1}^i - \eta].$$

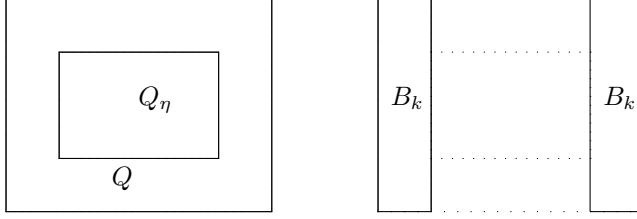
Now define the box,  $B^k$  by

$$B^k \equiv [\alpha_{j_1}^1, \alpha_{j_1+1}^1] \times \cdots \times [\alpha_{j_{k-1}}^{k-1}, \alpha_{j_{k-1}+1}^{k-1}] \times \\ [\alpha_{j_k}^k, \alpha_{j_k}^k + \eta] \times [\alpha_{j_{k+1}}^{k+1}, \alpha_{j_{k+1}+1}^{k+1}] \times \cdots \times [\alpha_{j_n}^n, \alpha_{j_n+1}^n]$$

or

$$B^k \equiv [\alpha_{j_1}^1, \alpha_{j_1+1}^1] \times \cdots \times [\alpha_{j_{k-1}}^{k-1}, \alpha_{j_{k-1}+1}^{k-1}] \times \\ [\alpha_{j_k}^k - \eta, \alpha_{j_k}^k] \times [\alpha_{j_{k+1}}^{k+1}, \alpha_{j_{k+1}+1}^{k+1}] \times \cdots \times [\alpha_{j_n}^n, \alpha_{j_n+1}^n].$$

In words, replace the closed interval in the  $k^{\text{th}}$  slot used to define  $Q$  with a much thinner closed interval at one end or the other while leaving the other intervals used to define  $Q$  the same. This is illustrated in the following picture.



The important thing to notice, is that every point of  $Q$  is either in  $Q_\eta$  or one of the sets,  $B_k$ . Therefore,

$$\begin{aligned}\mathcal{L}_{\mathcal{F}}(g) &\geq v(Q_\eta) - \sum_{k=1}^n 2Cv(B_k) \geq v(Q_\eta) - 4CL^{n-1}n\eta \\ &= v(Q_\eta) - K\eta\end{aligned}\tag{21.33}$$

where  $K$  is a constant which does not depend on  $\eta$ . Similarly,

$$\mathcal{U}_{\mathcal{F}}(g) \leq v(Q_\eta) + K\eta.\tag{21.34}$$

This implies  $\mathcal{U}_{\mathcal{F}}(g) - \mathcal{L}_{\mathcal{F}}(g) < 2K\eta$  and since  $\eta$  is arbitrary, the Riemann criterion verifies that  $g \in \mathcal{R}(\mathbb{R}^n)$ . Formulas (21.33) and (21.34) also verify that

$$\begin{aligned}v(Q_\eta) &\in [\mathcal{U}_{\mathcal{F}}(g) - K\eta, \mathcal{L}_{\mathcal{F}}(g) + K\eta] \\ &\subseteq [\mathcal{L}_{\mathcal{F}}(g) - K\eta, \mathcal{U}_{\mathcal{F}}(g) + K\eta].\end{aligned}$$

But also

$$\int g \, dx \in [\mathcal{L}_{\mathcal{F}}(g), \mathcal{U}_{\mathcal{F}}(g)] \subseteq [\mathcal{L}_{\mathcal{F}}(g) - K\eta, \mathcal{U}_{\mathcal{F}}(g) + K\eta]$$

and so

$$\left| \int g \, dx - v(Q_\eta) \right| \leq 4K\eta.$$

Now letting  $\eta \rightarrow 0$ , yields  $\int g \, dx = v(Q)$ .

Now let  $f$  be as described in the statement of the Corollary. Let  $f_Q$  be the value of  $f$  on  $\text{int}(Q)$ , and let  $g_Q$  be a function of the sort just considered which equals 1 on  $\text{int}(Q)$ . Then  $f$  is of the form

$$f = \sum_{Q \in \mathcal{G}} f_Q g_Q$$

with all but finitely many of the  $f_Q$  equal zero. Therefore, the above is really a finite sum and so by Corollary 21.11.10,  $f \in \mathcal{R}(\mathbb{R}^n)$  and

$$\int f \, dx = \sum_{Q \in \mathcal{G}} f_Q \int g_Q \, dx = \sum_{Q \in \mathcal{G}} f_Q v(Q).$$

There is a good deal of sloppiness inherent in the above description of a step function due to the fact that the boxes may be different but match up on an edge. It is convenient to be able to consider a more precise sort of function and this is done next.

For  $Q$  a box of the form

$$Q = \prod_{i=1}^k [a_i, b_i],$$

define the half open box,  $Q'$  by

$$Q' = \prod_{i=1}^k (a_i, b_i].$$

The reason for considering these sets is that if  $\mathcal{G}$  is a grid, the sets,  $Q'$  where  $Q \in \mathcal{G}$  are disjoint. Defining a step function,  $\phi$  as

$$\phi(\mathbf{x}) \equiv \sum_{Q \in \mathcal{G}} \phi_Q \chi_{Q'}(\mathbf{x}),$$

the number,  $\phi_Q$  is the value of  $\phi$  on the set,  $Q'$ . As before, define

$$M_{Q'}(f) \equiv \sup \{f(\mathbf{x}) : \mathbf{x} \in Q'\}, \quad m_{Q'}(f) \equiv \inf \{f(\mathbf{x}) : \mathbf{x} \in Q'\}.$$

The next lemma will be convenient a little later.

**Lemma 21.11.13** *Suppose  $f$  is a bounded function which equals zero off some bounded set. Then  $f \in \mathcal{R}(\mathbb{R}^n)$  if and only if for all  $\varepsilon > 0$  there exists a grid,  $\mathcal{G}$  such that*

$$\sum_{Q \in \mathcal{G}} (M_{Q'}(f) - m_{Q'}(f)) v(Q) < \varepsilon. \quad (21.35)$$

**Proof:** Since  $Q' \subseteq Q$ ,

$$M_{Q'}(f) - m_{Q'}(f) \leq M_Q(f) - m_Q(f)$$

and therefore, the only if part of the equivalence is obvious.

Conversely, let  $\mathcal{G}$  be a grid such that (21.35) holds with  $\varepsilon$  replaced with  $\frac{\varepsilon}{2}$ . It is necessary to show there is a grid such that (21.35) holds with no primes on the  $Q$ . Let  $\mathcal{F}$  be a refinement of  $\mathcal{G}$  obtained by adding the points  $\alpha_k^i + \eta_k$  where  $\eta_k \leq \eta$  and is also chosen so small that for each  $i = 1, \dots, n$ ,

$$\alpha_k^i + \eta_k < \alpha_{k+1}^i.$$

Then for

$$Q \equiv \prod_{i=1}^n [\alpha_{k_i}^i, \alpha_{k_i+1}^i] \in \mathcal{G},$$

Let

$$\hat{Q} \equiv \prod_{i=1}^n [\alpha_{k_i}^i + \eta_{k_i}, \alpha_{k_i+1}^i]$$

and denote by  $\hat{\mathcal{G}}$  the collection of these smaller boxes. For each set,  $Q$  in  $\mathcal{G}$  there is the smaller set,  $\hat{Q}$  along with  $n$  boxes,  $B_k, k = 1, \dots, n$ , one of whose sides is of length  $\eta_k$  and the remainder of whose sides are shorter than the diameter of  $Q$  such that the set,  $Q$  is the union of  $\hat{Q}$  and these sets,  $B_k$ . Now suppose  $f$  equals zero off the ball  $B(\mathbf{0}, \frac{R}{2})$ . Then without loss of generality, you may assume the diameter of every box in  $\mathcal{G}$  which has nonempty intersection with  $B(\mathbf{0}, R)$  is smaller than  $\frac{R}{3}$ . (If this is not so, simply refine  $\mathcal{G}$  to make it so, such a refinement leaving (21.35) valid.) Suppose there are  $P$  sets of  $\mathcal{G}$  contained in  $B(\mathbf{0}, R)$  and suppose that for all  $\mathbf{x}$ ,  $|f(\mathbf{x})| < C/2$ . Then

$$\sum_{Q \in \mathcal{F}} (M_Q(f) - m_Q(f)) v(Q) \leq \sum_{\hat{Q} \in \hat{\mathcal{G}}} (M_Q(f) - m_Q(f)) v(Q)$$

$$\begin{aligned}
& + \sum_{Q \in \mathcal{F} \setminus \hat{\mathcal{G}}} (M_Q(f) - m_Q(f)) v(Q) \\
& \leq \varepsilon/2 + CPnR^{n-1}\eta < \varepsilon
\end{aligned}$$

whenever  $\eta$  is small enough. Since  $\varepsilon$  is arbitrary,  $f \in \mathcal{R}(\mathbb{R}^n)$  as claimed.

**Definition 21.11.14** A bounded set,  $E$  is a Jordan set in  $\mathbb{R}^n$  or a contented set in  $\mathbb{R}^n$  if  $\chi_E \in \mathcal{R}(\mathbb{R}^n)$ . Also, for  $\mathcal{G}$  a grid and  $E$  a set, denote by  $\partial_{\mathcal{G}}(E)$  those boxes of  $\mathcal{G}$  which have nonempty intersection with both  $E$  and  $\mathbb{R}^n \setminus E$ .

The next theorem is a characterization of those sets which are Jordan sets.

**Theorem 21.11.15** A bounded set,  $E$ , is a Jordan set if and only if for every  $\varepsilon > 0$  there exists a grid,  $\mathcal{G}$ , such that

$$\sum_{Q \in \partial_{\mathcal{G}}(E)} v(Q) < \varepsilon.$$

**Proof:** If  $Q \notin \partial_{\mathcal{G}}(E)$ , then

$$M_Q(\chi_E) - m_Q(\chi_E) = 0$$

and if  $Q \in \partial_{\mathcal{G}}(E)$ , then

$$M_Q(\chi_E) - m_Q(\chi_E) = 1.$$

It follows that  $\mathcal{U}_{\mathcal{G}}(\chi_E) - \mathcal{L}_{\mathcal{G}}(\chi_E) = \sum_{Q \in \partial_{\mathcal{G}}(E)} v(Q)$  and this implies the conclusion of the theorem.

Note that if  $E$  is a Jordan set and if  $f \in \mathcal{R}(\mathbb{R}^n)$ , then by Corollary 21.11.10,  $\chi_E f \in \mathcal{R}(\mathbb{R}^n)$ .

**Definition 21.11.16** For  $E$  a Jordan set and  $f\chi_E \in \mathcal{R}(\mathbb{R}^n)$ .

$$\int_E f dV \equiv \int_{\mathbb{R}^n} \chi_E f dV.$$

A bounded set,  $E$ , has Jordan content 0 or content 0 if for every  $\varepsilon > 0$  there exists a grid,  $\mathcal{G}$  such that

$$\sum_{Q \cap E \neq \emptyset} v(Q) < \varepsilon.$$

This symbol says to sum the volumes of all boxes from  $\mathcal{G}$  which have nonempty intersection with  $E$ .

Note that any finite union of sets having Jordan content 0 also has Jordan content 0. (Why?)

**Definition 21.11.17** Let  $A$  be any subset of  $\mathbb{R}^n$ . Then  $\partial A$  denotes those points,  $\mathbf{x}$  with the property that if  $U$  is any open set containing  $\mathbf{x}$ , then  $U$  contains points of  $A$  as well as points of  $A^C$ .

**Corollary 21.11.18** If a bounded set,  $E \subseteq \mathbb{R}^n$  is contented, then  $\partial E$  has content 0.

**Proof:** Let  $\varepsilon > 0$  be given and suppose  $E$  is contented. Then there exists a grid,  $\mathcal{G}$  such that

$$\sum_{Q \in \partial_{\mathcal{G}}(E)} v(Q) < \frac{\varepsilon}{2n+1}. \quad (21.36)$$

Now refine  $\mathcal{G}$  if necessary to get a new grid,  $\mathcal{F}$  such that all boxes from  $\mathcal{F}$  which have nonempty intersection with  $\partial E$  have sides no larger than  $\delta$  where  $\delta$  is the smallest of all the sides of all the  $Q$  in the above sum. Recall that  $\partial_{\mathcal{G}}(E)$  consists of those boxes of  $\mathcal{G}$  which have nonempty intersection with both  $E$  and  $\mathbb{R}^n \setminus E$ .

Let  $\mathbf{x} \in \partial E$ . Then since the dimension is  $n$ , there are at most  $2^n$  boxes from  $\mathcal{F}$  which contain  $\mathbf{x}$ . Furthermore, at least one of these boxes is in  $\partial_{\mathcal{F}}(E)$  and is therefore a subset of a box from  $\partial_{\mathcal{G}}(E)$ . Here is why. If  $\mathbf{x}$  is an interior point of some  $Q \in \mathcal{F}$ , then there are points of both  $E$  and  $E^C$  contained in  $Q$  and so  $\mathbf{x} \in Q \in \partial_{\mathcal{F}}(E)$  and there are no other boxes from  $\mathcal{F}$  which contain  $\mathbf{x}$ . If  $\mathbf{x}$  is not an interior point of any  $Q \in \mathcal{F}$ , then the interior of the union of all the boxes from  $\mathcal{F}$  which do contain  $\mathbf{x}$  is an open set and therefore, must contain points of  $E$  and points from  $E^C$ . If  $\mathbf{x} \in E$ , then one of these boxes must contain points which are not in  $E$  since otherwise,  $\mathbf{x}$  would fail to be in  $\partial E$ . Pick that box. It is in  $\partial_{\mathcal{F}}(E)$  and contains  $\mathbf{x}$ . On the other hand, if  $\mathbf{x} \notin E$ , one of these boxes must contain points of  $E$  since otherwise,  $\mathbf{x}$  would fail to be in  $\partial E$ . Pick that box. This shows that every set from  $\mathcal{F}$  which contains a point of  $\partial E$  shares this point with a box of  $\partial_{\mathcal{G}}(E)$ . Let the boxes from  $\partial_{\mathcal{G}}(E)$  be  $\{P_1, \dots, P_m\}$ . Let  $\mathcal{S}(P_i)$  denote those sets of  $\mathcal{F}$  which contain a point of  $\partial E$  in common with  $P_i$ . Then if  $Q \in \mathcal{S}(P_i)$ , either  $Q \subseteq P_i$  or it intersects  $P_i$  on one of its  $2n$  faces. Therefore, the sum of the volumes of those boxes of  $\mathcal{S}(P_i)$  which intersect  $P_i$  on a particular face of  $P_i$  is no larger than  $v(P_i)$ . Consequently,

$$\sum_{Q \in \mathcal{S}(P_i)} v(Q) \leq 2nv(P_i) + v(P_i)$$

and so for  $Q \in \mathcal{F}$ ,

$$\sum_{Q \cap \partial E \neq \emptyset} v(Q) = \sum_{i=1}^m \sum_{Q \in \mathcal{S}(P_i)} v(Q) \leq \sum_{i=1}^m (2n+1)v(P_i) < \varepsilon$$

from (21.36). This proves the corollary.

**Theorem 21.11.19** *If a bounded set,  $E$ , has Jordan content 0, then  $E$  is a Jordan set and if  $f$  is any bounded function defined on  $E$ , then  $f\chi_E \in \mathcal{R}(\mathbb{R}^n)$  and*

$$\int_E f dV = 0.$$

**Proof:** Let  $\varepsilon > 0$ . Then let  $\mathcal{G}$  be a grid such that

$$\sum_{Q \cap E \neq \emptyset} v(Q) < \varepsilon.$$

Then every set of  $\partial_{\mathcal{G}}(E)$  contains a point of  $E$  so

$$\sum_{Q \in \partial_{\mathcal{G}}(E)} v(Q) \leq \sum_{Q \cap E \neq \emptyset} v(Q) < \varepsilon$$

and since  $\varepsilon$  was arbitrary, this shows from Theorem 21.11.15 that  $E$  is a Jordan set. Now let  $M$  be a positive number larger than all values of  $f$ , let  $m$  be a negative number smaller than all values of  $f$  and let  $\varepsilon > 0$  be given. Let  $\mathcal{G}$  be a grid with

$$\sum_{Q \cap E \neq \emptyset} v(Q) < \frac{\varepsilon}{1 + (M - m)}.$$

Then

$$\mathcal{U}_G(f\mathcal{X}_E) \leq \sum_{Q \cap E \neq \emptyset} Mv(Q) \leq \frac{\varepsilon M}{1 + (M - m)}$$

and

$$\mathcal{L}_G(f\mathcal{X}_E) \geq \sum_{Q \cap E \neq \emptyset} mv(Q) \geq \frac{\varepsilon m}{1 + (M - m)}$$

and so

$$\begin{aligned} \mathcal{U}_G(f\mathcal{X}_E) - \mathcal{L}_G(f\mathcal{X}_E) &\leq \sum_{Q \cap E \neq \emptyset} Mv(Q) - \sum_{Q \cap E \neq \emptyset} mv(Q) \\ &= (M - m) \sum_{Q \cap E \neq \emptyset} v(Q) < \frac{\varepsilon(m - N)}{1 + (M - m)} < \varepsilon. \end{aligned}$$

This shows  $f\mathcal{X}_E \in \mathcal{R}(\mathbb{R}^n)$ . Now also,

$$m\varepsilon \leq \int f\mathcal{X}_E dV \leq M\varepsilon$$

and since  $\varepsilon$  is arbitrary, this shows

$$\int_E f dV \equiv \int f\mathcal{X}_E dV = 0$$

and proves the theorem.

**Corollary 21.11.20** *If  $f\mathcal{X}_{E_i} \in \mathcal{R}(\mathbb{R}^n)$  for  $i = 1, 2, \dots, r$  and for all  $i \neq j$ ,  $E_i \cap E_j$  is either the empty set or a set of Jordan content 0, then letting  $F \equiv \cup_{i=1}^r E_i$ , it follows  $f\mathcal{X}_F \in \mathcal{R}(\mathbb{R}^n)$  and*

$$\int f\mathcal{X}_F dV \equiv \int_F f dV = \sum_{i=1}^r \int_{E_i} f dV.$$

**Proof:** By Corollary 21.11.10, this is true if  $r = 1$ . Suppose it is true for  $r$ . It will be shown that it is true for  $r + 1$ . Let  $F_r = \cup_{i=1}^r E_i$  and let  $F_{r+1}$  be defined similarly. By the induction hypothesis,  $f\mathcal{X}_{F_r} \in \mathcal{R}(\mathbb{R}^n)$ . Also, since  $F_r$  is a finite union of the  $E_i$ , it follows that  $F_r \cap E_{r+1}$  is either empty or a set of Jordan content 0.

$$-f\mathcal{X}_{F_r \cap E_{r+1}} + f\mathcal{X}_{F_r} + f\mathcal{X}_{E_{r+1}} = f\mathcal{X}_{F_{r+1}}$$

and by Theorem 21.11.19 each function on the left is in  $\mathcal{R}(\mathbb{R}^n)$  and the first one on the left has integral equal to zero. Therefore,

$$\int f\mathcal{X}_{F_{r+1}} dV = \int f\mathcal{X}_{F_r} dV + \int f\mathcal{X}_{E_{r+1}} dV$$

which by induction equals

$$\sum_{i=1}^r \int_{E_i} f dV + \int_{E_{r+1}} f dV = \sum_{i=1}^{r+1} \int_{E_i} f dV$$

and this proves the corollary.

What functions in addition to step functions are integrable? As in the case of integrals of functions of one variable, this is an important question. It turns out that the Riemann integrable functions are characterized by being continuous except on a very small set. To begin with it is necessary to define the oscillation of a function.



**Definition 21.11.21** Let  $f$  be a function defined on  $\mathbb{R}^n$  and let

$$\omega_{f,r}(\mathbf{x}) \equiv \sup \{|f(\mathbf{z}) - f(\mathbf{y})| : \mathbf{z}, \mathbf{y} \in B(\mathbf{x}, r)\}.$$

This is called the oscillation of  $f$  on  $B(\mathbf{x}, r)$ . Note that this function of  $r$  is decreasing in  $r$ . Define the oscillation of  $f$  as

$$\omega_f(\mathbf{x}) \equiv \lim_{r \rightarrow 0+} \omega_{f,r}(\mathbf{x}).$$

Note that as  $r$  decreases, the function,  $\omega_{f,r}(\mathbf{x})$  decreases. It is also bounded below by 0 and so the limit must exist and equals  $\inf \{\omega_{f,r}(\mathbf{x}) : r > 0\}$ . (Why?) Then the following simple lemma whose proof follows directly from the definition of continuity gives the reason for this definition.

**Lemma 21.11.22** A function,  $f$ , is continuous at  $\mathbf{x}$  if and only if  $\omega_f(\mathbf{x}) = 0$ .

This concept of oscillation gives a way to define how discontinuous a function is at a point. The discussion will depend on the following fundamental lemma which gives the existence of something called the Lebesgue number.

**Definition 21.11.23** Let  $\mathfrak{C}$  be a set whose elements are sets of  $\mathbb{R}^n$  and let  $K \subseteq \mathbb{R}^n$ . The set,  $\mathfrak{C}$  is called a cover of  $K$  if every point of  $K$  is contained in some set of  $\mathfrak{C}$ . If the elements of  $\mathfrak{C}$  are open sets, it is called an open cover.

**Lemma 21.11.24** Let  $K$  be sequentially compact and let  $\mathfrak{C}$  be an open cover of  $K$ . Then there exists  $r > 0$  such that whenever  $\mathbf{x} \in K$ ,  $B(\mathbf{x}, r)$  is contained in some set of  $\mathfrak{C}$ .

**Proof:** Suppose this is not so. Then letting  $r_n = 1/n$ , there exists  $\mathbf{x}_n \in K$  such that  $B(\mathbf{x}_n, r_n)$  is not contained in any set of  $\mathfrak{C}$ . Since  $K$  is sequentially compact, there is a subsequence,  $\mathbf{x}_{n_k}$  which converges to a point,  $\mathbf{x} \in K$ . But there exists  $\delta > 0$  such that  $B(\mathbf{x}, \delta) \subseteq U$  for some  $U \in \mathfrak{C}$ . Let  $k$  be so large that  $1/k < \delta/2$  and  $|\mathbf{x}_{n_k} - \mathbf{x}| < \delta/2$  also. Then if  $\mathbf{z} \in B(\mathbf{x}_{n_k}, r_{n_k})$ , it follows

$$|\mathbf{z} - \mathbf{x}| \leq |\mathbf{z} - \mathbf{x}_{n_k}| + |\mathbf{x}_{n_k} - \mathbf{x}| < \frac{\delta}{2} + \frac{\delta}{2} = \delta$$

and so  $B(\mathbf{x}_{n_k}, r_{n_k}) \subseteq U$  contrary to supposition. Therefore, the desired number exists after all.

**Theorem 21.11.25** Let  $f$  be a bounded function which equals zero off a bounded set and let  $W$  denote the set of points where  $f$  fails to be continuous. Then  $f \in \mathcal{R}(\mathbb{R}^n)$  if  $W$  has content zero. That is, for all  $\varepsilon > 0$  there exists a grid,  $\mathcal{G}$  such that

$$\sum_{Q \in \mathcal{G}_W} v(Q) < \varepsilon \quad (21.37)$$

where

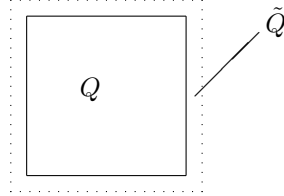
$$\mathcal{G}_W \equiv \{Q \in \mathcal{G} : Q \cap W \neq \emptyset\}.$$

**Proof:** Let  $|f(\mathbf{x})| < C/2$  for all  $\mathbf{x} \in \mathbb{R}^n$ , let  $\varepsilon > 0$  be given, and let  $\mathcal{G}$  be a grid which satisfies (21.37). Since  $f$  equals zero off some bounded set, there exists  $R$  such that  $f$  equals zero off of  $B(\mathbf{0}, \frac{R}{2})$ . Thus  $W \subseteq B(\mathbf{0}, \frac{R}{2})$ . Also note that if  $\mathcal{G}$  is a grid for which (21.37) holds, then this inequality continues to hold if  $\mathcal{G}$  is replaced with a refined grid. Therefore,

you may assume the diameter of every box in  $\mathcal{G}$  which intersects  $B(\mathbf{0}, R)$  is less than  $\frac{R}{3}$ . Since  $W$  is bounded,  $\mathcal{G}_W$  contains only finitely many boxes. Letting

$$Q \equiv \prod_{i=1}^n [a_i, b_i]$$

be one of these boxes, enlarge the box slightly as indicated in the following picture.



The enlarged box is an open set of the form,

$$\tilde{Q} \equiv \prod_{i=1}^n (a_i - \eta_i, b_i + \eta_i)$$

where  $\eta_i$  is chosen small enough that if

$$\prod_{i=1}^n (b_i + \eta_i - (a_i - \eta_i)) \equiv v(\tilde{Q}),$$

then

$$\sum_{Q \in \mathcal{G}_W} v(\tilde{Q}) < \varepsilon.$$

For each  $\mathbf{x} \in \mathbb{R}^n$ , let  $r_{\mathbf{x}}$  be such that

$$\omega_{f, r_{\mathbf{x}}}(\mathbf{x}) < \varepsilon + \omega_f(\mathbf{x}). \quad (21.38)$$

Now let  $\mathfrak{C}$  denote all intersections of the form  $\tilde{Q} \cap B(\mathbf{x}, r_{\mathbf{x}})$  such that  $\mathbf{x} \in \overline{B(\mathbf{0}, R)}$  so that  $\mathfrak{C}$  is an open cover of the compact set,  $\overline{B(\mathbf{0}, R)}$ . Let  $\delta$  be a Lebesgue number for this open cover of  $\overline{B(\mathbf{0}, R)}$  and let  $\mathcal{F}$  be a refinement of  $\mathcal{G}$  such that every box in  $\mathcal{F}$  has diameter less than  $\delta$ . Now let  $\mathcal{F}_1$  consist of those boxes of  $\mathcal{F}$  which have nonempty intersection with  $B(\mathbf{0}, R/2)$ . Thus all boxes of  $\mathcal{F}_1$  are contained in  $B(\mathbf{0}, R)$  and each one is contained in some set of  $\mathfrak{C}$ . Now let  $\mathfrak{C}_W$  be those open sets of  $\mathfrak{C}$ ,  $\tilde{Q} \cap B(\mathbf{x}, r_{\mathbf{x}})$ , for which  $\mathbf{x} \in W$  and let  $\mathcal{F}_W$  be those sets of  $\mathcal{F}_1$  which are subsets of some set of  $\mathfrak{C}_W$ . Then

$$\begin{aligned} \mathcal{U}_{\mathcal{F}}(f) - \mathcal{L}_{\mathcal{F}}(f) &= \sum_{Q \in \mathcal{F}_W} (M_Q(f) - m_Q(f)) v(Q) \\ &+ \sum_{Q \in \mathcal{F}_1 \setminus \mathcal{F}_W} (M_Q(f) - m_Q(f)) v(Q). \end{aligned}$$

If  $Q \in \mathcal{F}_1 \setminus \mathcal{F}_W$ , then  $Q$  must be a subset of some set of  $\mathfrak{C} \setminus \mathfrak{C}_W$  since it is not in any set of  $\mathfrak{C}_W$ . Therefore, from (21.38) and the observation that  $\mathbf{x} \notin W$ ,

$$M_Q(f) - m_Q(f) \leq \varepsilon.$$

Therefore,

$$\mathcal{U}_{\mathcal{F}}(f) - \mathcal{L}_{\mathcal{F}}(f) \leq \sum_{Q \in \mathcal{F}_W} C v(Q) + \sum_{Q \in \mathcal{F}_1 \setminus \mathcal{F}_W} \varepsilon v(Q)$$

$$\leq C\varepsilon + \varepsilon (2R)^n.$$

Since  $\varepsilon$  is arbitrary, this proves the theorem.<sup>4</sup>

From Theorem 21.11.15 you get a pretty good idea of what constitutes a contented set. These sets are essentially those which have thin boundaries. Most sets you are likely to think of will fall in this category. However, it is good to give specific examples of sets which are contented.

**Theorem 21.11.26** *Suppose  $E$  is a bounded contented set in  $\mathbb{R}^n$  and  $f, g : E \rightarrow \mathbb{R}$  are two functions satisfying  $f(\mathbf{x}) \geq g(\mathbf{x})$  for all  $\mathbf{x} \in E$  and  $f\chi_E$  and  $g\chi_E$  are both in  $\mathcal{R}(\mathbb{R}^n)$ . Now define*

$$P \equiv \{(\mathbf{x}, x_{n+1}) : \mathbf{x} \in E \text{ and } g(\mathbf{x}) \leq x_{n+1} \leq f(\mathbf{x})\}.$$

*Then  $P$  is a contented set in  $\mathbb{R}^{n+1}$ .*

**Proof:** Let  $\mathcal{G}$  be a grid such that for  $k = f, g$ ,

$$\mathcal{U}_{\mathcal{G}}(k) - \mathcal{L}_{\mathcal{G}}(k) < \varepsilon/2. \quad (21.39)$$

Let the boxes of  $\mathcal{G}$  which have nonempty intersection with  $E$  be  $\{Q_1, \dots, Q_m\}$  and let  $\{a_i\}_{i=-\infty}^{\infty}$  be a sequence on  $\mathbb{R}$ ,  $a_i < a_{i+1}$  for all  $i$ , which includes

$$M_{Q_j}(f\chi_E), M_{Q_j}(g\chi_E), m_{Q_j}(f\chi_E), m_{Q_j}(g\chi_E)$$

for all  $j = 1, \dots, m$ . Now define a grid on  $\mathbb{R}^{n+1}$  as follows.

$$\mathcal{G}' \equiv \{Q \times [a_i, a_{i+1}] : Q \in \mathcal{G}, i \in \mathbb{Z}\}$$

In words, this grid consists of all possible boxes of the form  $Q \times [a_i, a_{i+1}]$  where  $Q \in \mathcal{G}$  and  $a_i$  is a term of the sequence just described. It is necessary to verify that  $\chi_P \in \mathcal{R}(\mathbb{R}^{n+1})$ . This is done by showing that  $\mathcal{U}_{\mathcal{G}'}(\chi_P) - \mathcal{L}_{\mathcal{G}'}(\chi_P) < \varepsilon$  and then noting that  $\varepsilon > 0$  was arbitrary. For  $\mathcal{G}'$  just described, denote by  $Q'$  a box in  $\mathcal{G}'$ . Thus  $Q' = Q \times [a_i, a_{i+1}]$  for some  $i$ .

$$\begin{aligned} \mathcal{U}_{\mathcal{G}'}(\chi_P) - \mathcal{L}_{\mathcal{G}'}(\chi_P) &\equiv \sum_{Q' \in \mathcal{G}'} (M_{Q'}(\chi_P) - m_{Q'}(\chi_P)) v_{n+1}(Q') \\ &= \sum_{i=-\infty}^{\infty} \sum_{j=1}^m \left( M_{Q_j}(\chi_P) - m_{Q_j}(\chi_P) \right) v_n(Q_j) (a_{i+1} - a_i) \end{aligned}$$

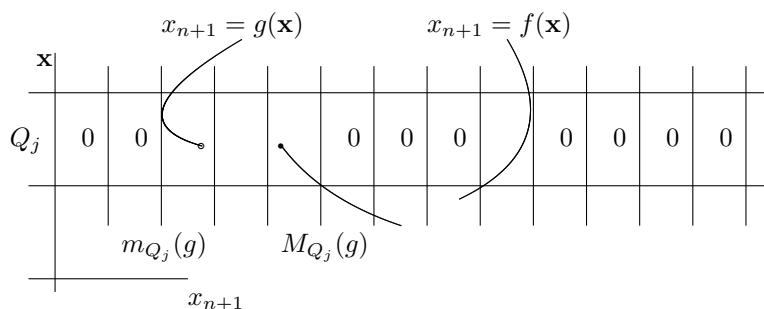
and all sums are bounded because the functions,  $f$  and  $g$  are given to be bounded. Therefore, there are no limit considerations needed here. Thus

$$\mathcal{U}_{\mathcal{G}'}(\chi_P) - \mathcal{L}_{\mathcal{G}'}(\chi_P) =$$

$$\sum_{j=1}^m v_n(Q_j) \sum_{i=-\infty}^{\infty} (M_{Q_j \times [a_i, a_{i+1}]}(\chi_P) - m_{Q_j \times [a_i, a_{i+1}]}(\chi_P)) (a_{i+1} - a_i).$$

Consider the inside sum with the aid of the following picture.

<sup>4</sup>In fact one cannot do any better. It can be shown that if a function is Riemann integrable, then it must be the case that for all  $\varepsilon > 0$ , (21.37) is satisfied for some grid,  $\mathcal{G}$ . This along with what was just shown is known as Lebesgue's theorem after Lebesgue who discovered it in the early years of the twentieth century. Actually, he also invented a far superior integral which has been the integral of serious mathematicians since that time.



In this picture, the little rectangles represent the boxes  $Q_j \times [a_i, a_{i+1}]$  for fixed  $j$ . The part of  $P$  having  $\mathbf{x}$  contained in  $Q_j$  is between the two surfaces,  $x_{n+1} = g(\mathbf{x})$  and  $x_{n+1} = f(\mathbf{x})$  and there is a zero placed in those boxes for which  $M_{Q_j \times [a_i, a_{i+1}]}(\mathcal{X}_P) - m_{Q_j \times [a_i, a_{i+1}]}(\mathcal{X}_P) = 0$ . You see,  $\mathcal{X}_P$  has either the value of 1 or the value of 0 depending on whether  $(\mathbf{x}, y)$  is contained in  $P$ . For the boxes shown with 0 in them, either all of the box is contained in  $P$  or none of the box is contained in  $P$ . Either way,  $M_{Q_j \times [a_i, a_{i+1}]}(\mathcal{X}_P) - m_{Q_j \times [a_i, a_{i+1}]}(\mathcal{X}_P) = 0$  on these boxes. However, on the boxes intersected by the surfaces, the value of  $M_{Q_j \times [a_i, a_{i+1}]}(\mathcal{X}_P) - m_{Q_j \times [a_i, a_{i+1}]}(\mathcal{X}_P)$  is 1 because there are points in this box which are not in  $P$  as well as points which are in  $P$ . Because of the construction of  $\mathcal{G}'$  which included all values of  $M_{Q_j}(f\mathcal{X}_E), M_{Q_j}(g\mathcal{X}_E), m_{Q_j}(f\mathcal{X}_E), m_{Q_j}(g\mathcal{X}_E)$  for all  $j = 1, \dots, m$ ,

$$\begin{aligned} & \sum_{i=-\infty}^{\infty} (M_{Q_j \times [a_i, a_{i+1}]}(\mathcal{X}_P) - m_{Q_j \times [a_i, a_{i+1}]}(\mathcal{X}_P)) (a_{i+1} - a_i) \leq \\ & \sum_{\{i: m_{Q_j}(g) \leq a_i < M_{Q_j}(g)\}} 1 (a_{i+1} - a_i) + \sum_{\{i: m_{Q_j}(f) \leq a_i < M_{Q_j}(f)\}} 1 (a_{i+1} - a_i) \\ & = (M_{Q_j}(g) - m_{Q_j}(g)) + (M_{Q_j}(f) - m_{Q_j}(f)). \end{aligned}$$

(Note the inequality.) Therefore, by (21.39),

$$\begin{aligned} \mathcal{U}_{\mathcal{G}'}(\mathcal{X}_P) - \mathcal{L}_{\mathcal{G}'}(\mathcal{X}_P) & \leq \sum_{j=1}^m v_n(Q_j) [(M_{Q_j}(g) - m_{Q_j}(g)) + (M_{Q_j}(f) - m_{Q_j}(f))] \\ & = \mathcal{U}_{\mathcal{G}}(f) - \mathcal{L}_{\mathcal{G}}(f) + \mathcal{U}_{\mathcal{G}}(g) - \mathcal{L}_{\mathcal{G}}(g) < \varepsilon/2 + \varepsilon/2 = \varepsilon. \end{aligned}$$

Since  $\varepsilon > 0$  is arbitrary, this proves the theorem.

**Corollary 21.11.27** Suppose  $f$  and  $g$  are continuous functions defined on  $E$ , a contented set in  $\mathbb{R}^n$  and that  $g(\mathbf{x}) \leq f(\mathbf{x})$  for all  $\mathbf{x} \in E$ . Then

$$P \equiv \{(\mathbf{x}, x_{n+1}) : \mathbf{x} \in E \text{ and } g(\mathbf{x}) \leq x_{n+1} \leq f(\mathbf{x})\}$$

is a contented set in  $\mathbb{R}^n$ .

**Proof:** Extend  $f$  and  $g$  to equal 0 off  $E$ . The set of discontinuities of  $f$  and  $g$  is contained in  $\partial E$  and Corollary 21.11.18 on Page 566 implies this is a set of content 0. Therefore, from Theorem 21.11.25, for  $k = f, g$ , it follows that  $k\mathcal{X}_E$  is in  $\mathcal{R}(\mathbb{R}^n)$  because the set of discontinuities is contained in  $\partial E$ . The conclusion now follows from Theorem 21.11.26. This proves the corollary.

As an example of how this can be applied, it is obvious a closed interval is a contented set in  $\mathbb{R}$ . Therefore, if  $f, g$  are two continuous functions with  $f(x) \geq g(x)$  for  $x \in [a, b]$ , it follows from the above theorem or its corollary that the set,

$$P_1 \equiv \{(x, y) : g(x) \leq y \leq f(x)\}$$

is a contented set in  $\mathbb{R}^2$ . Now using the theorem and corollary again, suppose  $f_1(x, y) \geq g_1(x, y)$  for  $(x, y) \in P_1$  and  $f, g$  are continuous. Then the set

$$P_2 \equiv \{(x, y, z) : g_1(x, y) \leq z \leq f_1(x, y)\}$$

is a contented set in  $\mathbb{R}^3$ . Clearly you can continue this way obtaining examples of contented sets.

Note that as a special case of Corollary 21.11.12 on Page 563, it follows that every box is a contented set.

### 21.11.2 Iterated Integrals

To evaluate an  $n$  dimensional Riemann integral, one uses iterated integrals. Formally, an iterated integral is defined as follows. For  $f$  a function defined on  $\mathbb{R}^{n+m}$ ,

$$\mathbf{y} \rightarrow f(\mathbf{x}, \mathbf{y})$$

is a function of  $\mathbf{y}$  for each  $\mathbf{x} \in \mathbb{R}^{n+m}$ . Therefore, it might be possible to integrate this function of  $\mathbf{y}$  and write

$$\int_{\mathbb{R}^m} f(\mathbf{x}, \mathbf{y}) dV_{\mathbf{y}}.$$

Now the result is clearly a function of  $\mathbf{x}$  and so, it might be possible to integrate this and write

$$\int_{\mathbb{R}^n} \int_{\mathbb{R}^m} f(\mathbf{x}, \mathbf{y}) dV_{\mathbf{y}} dV_{\mathbf{x}}.$$

This symbol is called an iterated integral, because it involves the iteration of two lower dimensional integrations. Under what conditions are the two iterated integrals equal to the integral

$$\int_{\mathbb{R}^{n+m}} f(\mathbf{z}) dV?$$

**Definition 21.11.28** Let  $\mathcal{G}$  be a grid on  $\mathbb{R}^{n+m}$  defined by the  $n+m$  sequences,

$$\{\alpha_k^i\}_{k=-\infty}^{\infty} \quad i = 1, \dots, n+m.$$

Let  $\mathcal{G}_n$  be the grid on  $\mathbb{R}^n$  obtained by considering only the first  $n$  of these sequences and let  $\mathcal{G}_m$  be the grid on  $\mathbb{R}^m$  obtained by considering only the last  $m$  of the sequences. Thus a typical box in  $\mathcal{G}_m$  would be

$$\prod_{i=n+1}^{n+m} [\alpha_{k_i}^i, \alpha_{k_i+1}^i], \quad k_i \geq n+1$$

and a box in  $\mathcal{G}_n$  would be of the form

$$\prod_{i=1}^n [\alpha_{k_i}^i, \alpha_{k_i+1}^i], \quad k_i \leq n.$$

**Lemma 21.11.29** *Let  $\mathcal{G}$ ,  $\mathcal{G}_n$ , and  $\mathcal{G}_m$  be the grids defined above. Then*

$$\mathcal{G} = \{R \times P : R \in \mathcal{G}_n \text{ and } P \in \mathcal{G}_m\}.$$

**Proof:** If  $Q \in \mathcal{G}$ , then  $Q$  is clearly of this form. On the other hand, if  $R \times P$  is one of the sets described above, then from the above description of  $R$  and  $P$ , it follows  $R \times P$  is one of the sets of  $\mathcal{G}$ .

Now let  $\mathcal{G}$  be a grid on  $\mathbb{R}^{n+m}$  and suppose

$$\phi(\mathbf{z}) = \sum_{Q \in \mathcal{G}} \phi_Q \mathcal{X}_{Q'}(\mathbf{z}) \quad (21.40)$$

where  $\phi_Q$  equals zero for all but finitely many  $Q$ . Thus  $\phi$  is a step function. Recall that for

$$Q = \prod_{i=1}^{n+m} [a_i, b_i], \quad Q' \equiv \prod_{i=1}^{n+m} (a_i, b_i]$$

Letting  $(\mathbf{x}, \mathbf{y}) = \mathbf{z}$ , Lemma 21.11.29 implies

$$\begin{aligned} \phi(\mathbf{z}) &= \phi(\mathbf{x}, \mathbf{y}) = \sum_{R \in \mathcal{G}_n} \sum_{P \in \mathcal{G}_m} \phi_{R \times P} \mathcal{X}_{R' \times P'}(\mathbf{x}, \mathbf{y}) \\ &= \sum_{R \in \mathcal{G}_n} \sum_{P \in \mathcal{G}_m} \phi_{R \times P} \mathcal{X}_{R'}(\mathbf{x}) \mathcal{X}_{P'}(\mathbf{y}). \end{aligned} \quad (21.41)$$

For a function of two variables,  $h$ , denote by  $h(\cdot, \mathbf{y})$  the function,  $\mathbf{x} \rightarrow h(\mathbf{x}, \mathbf{y})$  and  $h(\mathbf{x}, \cdot)$  the function  $\mathbf{y} \rightarrow h(\mathbf{x}, \mathbf{y})$ . The following lemma is a preliminary version of Fubini's theorem.

**Lemma 21.11.30** *Let  $\phi$  be a step function as described in (21.40). Then*

$$\phi(\mathbf{x}, \cdot) \in \mathcal{R}(\mathbb{R}^m), \quad (21.42)$$

$$\int_{\mathbb{R}^m} \phi(\cdot, \mathbf{y}) dV_y \in \mathcal{R}(\mathbb{R}^n), \quad (21.43)$$

and

$$\int_{\mathbb{R}^n} \int_{\mathbb{R}^m} \phi(\mathbf{x}, \mathbf{y}) dV_y dV_x = \int_{\mathbb{R}^{n+m}} \phi(\mathbf{z}) dV. \quad (21.44)$$

**Proof:** To verify (21.42), note that  $\phi(\mathbf{x}, \cdot)$  is the step function

$$\phi(\mathbf{x}, \mathbf{y}) = \sum_{P \in \mathcal{G}_m} \phi_{R \times P} \mathcal{X}_{P'}(\mathbf{y}).$$

Where  $\mathbf{x} \in R'$ . By Corollary 21.11.12, this verifies (21.42). From the description in (21.41) and this corollary,

$$\begin{aligned} \int_{\mathbb{R}^m} \phi(\mathbf{x}, \mathbf{y}) dV_y &= \sum_{R \in \mathcal{G}_n} \sum_{P \in \mathcal{G}_m} \phi_{R \times P} \mathcal{X}_{R'}(\mathbf{x}) v(P) \\ &= \sum_{R \in \mathcal{G}_n} \left( \sum_{P \in \mathcal{G}_m} \phi_{R \times P} v(P) \right) \mathcal{X}_{R'}(\mathbf{x}), \end{aligned} \quad (21.45)$$

another step function. Therefore, Corollary 21.11.12 applies again to verify (21.43). Finally, (21.45) implies

$$\begin{aligned} \int_{\mathbb{R}^n} \int_{\mathbb{R}^m} \phi(\mathbf{x}, \mathbf{y}) \, dV_y \, dV_x &= \sum_{R \in \mathcal{G}_n} \sum_{P \in \mathcal{G}_m} \phi_{R \times P} v(P) v(R) \\ &= \sum_{Q \in \mathcal{G}} \phi_Q v(Q) = \int_{\mathbb{R}^{n+m}} \phi(\mathbf{z}) \, dV. \end{aligned}$$

and this proves the lemma.

From (21.45),

$$\begin{aligned} M_{R'_1} \left( \int_{\mathbb{R}^m} \phi(\cdot, \mathbf{y}) \, dV_y \right) &\equiv \sup \left\{ \sum_{R \in \mathcal{G}_n} \left( \sum_{P \in \mathcal{G}_m} \phi_{R \times P} v(P) \right) \chi_{R'}(\mathbf{x}) : \mathbf{x} \in R'_1 \right\} \\ &= \sum_{P \in \mathcal{G}_m} \phi_{R_1 \times P} v(P). \end{aligned} \quad (21.46)$$

Similarly,

$$\begin{aligned} m_{R'_1} \left( \int_{\mathbb{R}^m} \phi(\cdot, \mathbf{y}) \, dV_y \right) &\equiv \inf \left\{ \sum_{R \in \mathcal{G}_n} \left( \sum_{P \in \mathcal{G}_m} \phi_{R \times P} v(P) \right) \chi_{R'}(\mathbf{x}) : \mathbf{x} \in R'_1 \right\} \\ &= \sum_{P \in \mathcal{G}_m} \phi_{R_1 \times P} v(P). \end{aligned} \quad (21.47)$$

**Theorem 21.11.31 (Fubini)** Let  $f \in \mathcal{R}(\mathbb{R}^{n+m})$  and suppose also that  $f(\mathbf{x}, \cdot) \in \mathcal{R}(\mathbb{R}^m)$  for each  $\mathbf{x}$ . Then

$$\int_{\mathbb{R}^m} f(\cdot, \mathbf{y}) \, dV_y \in \mathcal{R}(\mathbb{R}^n) \quad (21.48)$$

and

$$\int_{\mathbb{R}^{n+m}} f(\mathbf{z}) \, dV = \int_{\mathbb{R}^n} \int_{\mathbb{R}^m} f(\mathbf{x}, \mathbf{y}) \, dV_y \, dV_x. \quad (21.49)$$

**Proof:** Let  $\mathcal{G}$  be a grid such that  $\mathcal{U}_{\mathcal{G}}(f) - \mathcal{L}_{\mathcal{G}}(f) < \varepsilon$  and let  $\mathcal{G}_n$  and  $\mathcal{G}_m$  be as defined above. Let

$$\phi(\mathbf{z}) \equiv \sum_{Q \in \mathcal{G}} M_{Q'}(f) \chi_{Q'}(\mathbf{z}), \quad \psi(\mathbf{z}) \equiv \sum_{Q \in \mathcal{G}} m_{Q'}(f) \chi_{Q'}(\mathbf{z}).$$

By Corollary 21.11.12, and the observation that  $M_{Q'}(f) \leq M_Q(f)$  and  $m_{Q'}(f) \geq m_Q(f)$ ,

$$\mathcal{U}_{\mathcal{G}}(f) \geq \int \phi \, dV, \quad \mathcal{L}_{\mathcal{G}}(f) \leq \int \psi \, dV.$$

Also  $f(\mathbf{z}) \in (\psi(\mathbf{z}), \phi(\mathbf{z}))$  for all  $\mathbf{z}$ . Thus from (21.46),

$$M_{R'} \left( \int_{\mathbb{R}^m} f(\cdot, \mathbf{y}) \, dV_y \right) \leq M_{R'} \left( \int_{\mathbb{R}^m} \phi(\cdot, \mathbf{y}) \, dV_y \right) = \sum_{P \in \mathcal{G}_m} M_{R' \times P'}(f) v(P)$$

and from (21.47),

$$m_{R'} \left( \int_{\mathbb{R}^m} f(\cdot, \mathbf{y}) \, dV_y \right) \geq m_{R'} \left( \int_{\mathbb{R}^m} \psi(\cdot, \mathbf{y}) \, dV_y \right) = \sum_{P \in \mathcal{G}_m} m_{R' \times P'}(f) v(P).$$

Therefore,

$$\sum_{R \in \mathcal{G}_n} \left[ M_{R'} \left( \int_{\mathbb{R}^m} f(\cdot, \mathbf{y}) dV_y \right) - m_{R'} \left( \int_{\mathbb{R}^m} f(\cdot, \mathbf{y}) dV_y \right) \right] v(R) \leq$$

$$\sum_{R \in \mathcal{G}_n} \sum_{P \in \mathcal{G}_m} [M_{R' \times P'}(f) - m_{R' \times P'}(f)] v(P) v(R) \leq \mathcal{U}_{\mathcal{G}}(f) - \mathcal{L}_{\mathcal{G}}(f) < \varepsilon.$$

This shows, from Lemma 21.11.13, that  $\int_{\mathbb{R}^m} f(\cdot, \mathbf{y}) dV_y \in \mathcal{R}(\mathbb{R}^n)$ . It remains to verify (21.49). First note

$$\int_{\mathbb{R}^{n+m}} f(\mathbf{z}) dV \in [\mathcal{L}_{\mathcal{G}}(f), \mathcal{U}_{\mathcal{G}}(f)].$$

Next, by Lemma 21.11.30,

$$\mathcal{L}_{\mathcal{G}}(f) \leq \int_{\mathbb{R}^{n+m}} \psi dV = \int_{\mathbb{R}^n} \int_{\mathbb{R}^m} \psi dV_y dV_x \leq \int_{\mathbb{R}^n} \int_{\mathbb{R}^m} f(\mathbf{x}, \mathbf{y}) dV_y dV_x$$

$$\leq \int_{\mathbb{R}^n} \int_{\mathbb{R}^m} \phi(\mathbf{x}, \mathbf{y}) dV_y dV_x = \int_{\mathbb{R}^{n+m}} \phi dV \leq \mathcal{U}_{\mathcal{G}}(f).$$

Therefore,

$$\left| \int_{\mathbb{R}^n} \int_{\mathbb{R}^m} f(\mathbf{x}, \mathbf{y}) dV_y dV_x - \int_{\mathbb{R}^{n+m}} f(\mathbf{z}) dV \right| \leq \varepsilon$$

and since  $\varepsilon > 0$  is arbitrary, this proves Fubini's theorem<sup>5</sup>.

**Corollary 21.11.32** *Suppose  $E$  is a bounded contented set in  $\mathbb{R}^n$  and let  $\phi, \psi$  be continuous functions defined on  $E$  such that  $\phi(\mathbf{x}) \geq \psi(\mathbf{x})$ . Also suppose  $f$  is a continuous bounded function defined on the set,*

$$P \equiv \{(\mathbf{x}, y) : \psi(\mathbf{x}) \leq y \leq \phi(\mathbf{x})\},$$

*It follows  $f\mathcal{X}_P \in \mathcal{R}(\mathbb{R}^{n+1})$  and*

$$\int_P f dV = \int_E \int_{\psi(\mathbf{x})}^{\phi(\mathbf{x})} f(\mathbf{x}, y) dy dV_x.$$

**Proof:** Since  $f$  is continuous, there is no problem in writing  $f(\mathbf{x}, \cdot) \in \mathcal{R}(\mathbb{R}^1)$ . Also,  $f\mathcal{X}_P \in \mathcal{R}(\mathbb{R}^{n+1})$  because  $P$  is contented thanks to Corollary 21.11.27. Therefore, by Fubini's theorem

$$\begin{aligned} \int_P f dV &= \int_{\mathbb{R}^n} \int_{\mathbb{R}} f\mathcal{X}_P dy dV_x \\ &= \int_E \int_{\psi(\mathbf{x})}^{\phi(\mathbf{x})} f(\mathbf{x}, y) dy dV_x \end{aligned}$$

proving the corollary.

Other versions of this corollary are immediate and should be obvious whenever encountered.

---

<sup>5</sup>Actually, Fubini's theorem usually refers to a much more profound result in the theory of Lebesgue integration.



### 21.11.3 Some Observations

Some of the above material is very technical. This is because it gives complete answers to the fundamental questions on existence of the integral and related theoretical considerations. It was realized early in the twentieth century that these difficulties occur because from the point of view of mathematics, this is not the right way to define an integral! Better results are obtained much more easily using the Lebesgue integral. Many of the technicalities related to Jordan content disappear almost magically when the right integral is used. However, the Lebesgue integral is much more abstract than the Riemann integral and it is not traditional to consider it in a beginning calculus course. If you are interested in the fundamental properties of the integral and the theory behind it, you should abandon the Riemann integral which is an antiquated relic and begin to study the integral of the last century. One of the best introductions to it is in [14]. Another very good source is [8]. This advanced calculus text does everything in terms of the Lebesgue integral and never bothers to struggle with the inferior Riemann integral. A more general treatment is found in [10], [11], [15], and [12]. There is also a still more general integral called the generalized Riemann integral. A recent book on this subject is [3].



# The Integral On Other Sets

Consider the boundary of some three dimensional region such that a function,  $f$  is defined on this boundary. Imagine taking the value of this function at a point, multiplying this value by the area of an infinitesimal chunk of area located at this point and then adding these up. This is just the notion of the integral presented earlier only now there is a difference because this infinitesimal chunk of area should be considered as two dimensional even though it is in three dimensions. However, it is not really all that different from what was done earlier. As before, it all depends on the following fundamental definition on Page 460.

**Definition 22.0.33** Let  $\mathbf{u}_1, \dots, \mathbf{u}_p$  be vectors in  $\mathbb{R}^n$ . The  $p$  dimensional parallelepiped determined by these vectors will be denoted by  $P(\mathbf{u}_1, \dots, \mathbf{u}_p)$  and it is defined as

$$P(\mathbf{u}_1, \dots, \mathbf{u}_p) \equiv \left\{ \sum_{j=1}^p s_j \mathbf{u}_j : s_j \in [0, 1] \right\}.$$

Define the volume of this parallelepiped by

$$\text{volume of } P(\mathbf{u}_1, \dots, \mathbf{u}_p) \equiv (\det(\mathbf{u}_i \cdot \mathbf{u}_j))^{1/2}.$$

Suppose then that  $\mathbf{x} = \mathbf{f}(\mathbf{u})$  where  $\mathbf{u} \in U$ , a subset of  $\mathbb{R}^p$  and  $\mathbf{x}$  is a point in  $V$ , a subset of  $n$  dimensional space where  $n \geq p$ . Thus, letting the Cartesian coordinates of  $\mathbf{x}$  be given by  $\mathbf{x} = (x_1, \dots, x_n)^T$ , each  $x_i$  being a function of  $\mathbf{u}$ , an infinitesimal box located at  $\mathbf{u}_0$  corresponds to an infinitesimal parallelepiped located at  $\mathbf{f}(\mathbf{u}_0)$  which is determined by the  $p$  vectors  $\left\{ \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_i} du_i \right\}_{i=1}^p$ , each of which is tangent to the surface defined by  $\mathbf{x} = \mathbf{f}(\mathbf{u})$ . (No sum on the repeated index.) From Definition 22.0.33, the volume of this infinitesimal parallelepiped located at  $\mathbf{f}(\mathbf{u}_0)$  is given by

$$\det \left( \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_i} du_i \cdot \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_j} du_j \right)^{1/2}. \quad (22.1)$$

I like to think of this in the case where  $p = 2$ . In this case the infinitesimal parallelepiped is an infinitesimal parallelogram tangent to the surface defined by  $\mathbf{x} = \mathbf{f}(\mathbf{u})$  like a very small scale on a lizard. This is the essence of the idea. To define the area of the lizard sum up areas of individual scales.

Now, continuing with the general case, the matrix in the above formula is a  $p \times p$  matrix. Denoting

$$\frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_i} = \mathbf{x}_{,i}$$

to save space, this matrix is of the form

$$\begin{pmatrix} \overbrace{\begin{pmatrix} du_1 & 0 & \cdots & 0 \\ 0 & du_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & du_p \end{pmatrix}}^{p \times p} \overbrace{\begin{pmatrix} \mathbf{x}_{,1}^T \\ \mathbf{x}_{,2}^T \\ \vdots \\ \mathbf{x}_{,p}^T \end{pmatrix}}^{p \times n} \\ \left( \overbrace{\begin{pmatrix} \mathbf{x}_{,1} & \mathbf{x}_{,2} & \cdots & \mathbf{x}_{,p} \end{pmatrix}}^{n \times p} \right) \overbrace{\begin{pmatrix} du_1 & 0 & \cdots & 0 \\ 0 & du_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & du_p \end{pmatrix}}^{p \times p} \end{pmatrix}$$

Therefore, by the theorem which says the determinant of a product equals the product of the determinants, Theorem 19.9.11 on Page 452, the determinant of the above product equals

$$\det \begin{pmatrix} du_1 & 0 & \cdots & 0 \\ 0 & du_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & du_p \end{pmatrix}^2 \det \left( \begin{pmatrix} \mathbf{x}_{,1}^T \\ \mathbf{x}_{,2}^T \\ \vdots \\ \mathbf{x}_{,p}^T \end{pmatrix} \begin{pmatrix} \mathbf{x}_{,1} & \mathbf{x}_{,2} & \cdots & \mathbf{x}_{,p} \end{pmatrix} \right) =$$

$$\det \left( \begin{pmatrix} \mathbf{x}_{,1}^T \\ \mathbf{x}_{,2}^T \\ \vdots \\ \mathbf{x}_{,p}^T \end{pmatrix} \begin{pmatrix} \mathbf{x}_{,1} & \mathbf{x}_{,2} & \cdots & \mathbf{x}_{,p} \end{pmatrix} \right) (du_1 du_2 \cdots du_p)^2$$

and so taking the square root implies the volume of the infinitesimal parallelepiped at  $\mathbf{x} = \mathbf{f}(\mathbf{u}_0)$  is

$$\det \left( \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_i} \cdot \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_j} \right)^{1/2} du_1 du_2 \cdots du_p =$$

$$\det \left( D\mathbf{f}(\mathbf{u})^T D\mathbf{f}(\mathbf{u}) \right)^{1/2} du_1 du_2 \cdots du_p$$

**Definition 22.0.34** Let  $\mathbf{x} = \mathbf{f}(\mathbf{u})$  be as described above. Then the symbol,  $\frac{\partial(x_1, \dots, x_n)}{\partial(u_1, \dots, u_p)}$ , is defined by

$$\det \left( \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_i} \cdot \frac{\partial \mathbf{x}(\mathbf{u}_0)}{\partial u_j} \right)^{1/2} \equiv \frac{\partial(x_1, \dots, x_n)}{\partial(u_1, \dots, u_p)}.$$

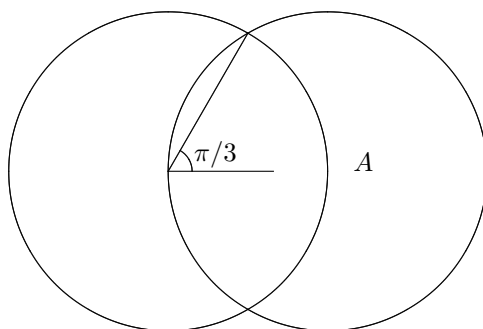
Also, the symbol,  $dV_p \equiv \frac{\partial(x_1, \dots, x_n)}{\partial(u_1, \dots, u_p)} du_1 \cdots du_p$  is called the volume element or area element. Note the use of the subscript,  $p$ . This indicates the  $p$  dimensional volume element. When  $p = 2$  it is customary to write  $dA$ . Also, continue referring to  $\frac{\partial(x_1, \dots, x_n)}{\partial(u_1, \dots, u_p)}$  as the Jacobian.

This motivates the following fundamental procedure which I hope is extremely familiar from the earlier material.

**Procedure 22.0.35** Suppose  $U$  is a subset of  $\mathbb{R}^p$  and suppose  $\mathbf{f} : U \rightarrow \mathbf{f}(U) \subseteq \mathbb{R}^n$  is a one to one and  $C^1$  function. Then if  $h : \mathbf{f}(U) \rightarrow \mathbb{R}$ , define the  $p$  dimensional surface integral,  $\int_{\mathbf{f}(U)} h(\mathbf{x}) dV_p$  according to the following formula.

$$\int_{\mathbf{f}(U)} h(\mathbf{x}) dV_p \equiv \int_U h(\mathbf{f}(\mathbf{u})) \frac{\partial(x_1, \dots, x_n)}{\partial(u_1, \dots, u_p)} dV.$$

**Example 22.0.36** Find the area of the region labeled  $A$  in the following picture. The two circles are of radius 1, one has center  $(0,0)$  and the other has center  $(1,0)$ .



The circles bounding these disks are  $x^2 + y^2 = 1$  and  $(x-1)^2 + y^2 = x^2 + y^2 - 2x + 1 = 1$ . Therefore, in polar coordinates these are of the form  $r = 1$  and  $r = 2 \cos \theta$ .

The set  $A$  corresponds to the set  $U$ , in the  $(\theta, r)$  plane determined by  $\theta \in [-\frac{\pi}{3}, \frac{\pi}{3}]$  and for each value of  $\theta$  in this interval,  $r$  goes from 1 up to  $2 \cos \theta$ . Therefore, the area of this region is of the form,

$$\int_U 1 dV = \int_{-\pi/3}^{\pi/3} \int_1^{2 \cos \theta} \frac{\partial(x_1, x_2)}{\partial(\theta, r)} dr d\theta.$$

It is necessary to find  $\frac{\partial(x_1, x_2)}{\partial(\theta, r)}$ . The mapping  $\mathbf{f} : U \rightarrow \mathbb{R}^2$  takes the form  $\mathbf{f}(\theta, r) = (r \cos \theta, r \sin \theta)^T$  and so

$$D\mathbf{f}(\theta, r) = \begin{pmatrix} -r \sin \theta & \cos \theta \\ r \cos \theta & \sin \theta \end{pmatrix}$$

and so

$$D\mathbf{f}(\theta, r)^T D\mathbf{f}(\theta, r) = \begin{pmatrix} r^2 & 0 \\ 0 & 1 \end{pmatrix}$$

which implies

$$\frac{\partial(x_1, x_2)}{\partial(\theta, r)} = \det \left( D\mathbf{f}(\theta, r)^T D\mathbf{f}(\theta, r) \right)^{1/2} = r.$$

Therefore, the area element is  $r dr d\theta$ . It follows the desired area is

$$\int_{-\pi/3}^{\pi/3} \int_1^{2 \cos \theta} r dr d\theta = \frac{1}{2} \sqrt{3} + \frac{1}{3} \pi.$$

**Example 22.0.37** Consider the surface given by  $z = x^2$  for  $(x, y) \in [0, 1] \times [0, 1] = U$ . Find the surface area of this surface.

The first step in using the above is to write this surface in the form  $\mathbf{x} = \mathbf{f}(\mathbf{u})$ . This is easy to do if you let  $\mathbf{u} = (x, y)$ . Then  $\mathbf{f}(x, y) = (x, y, x^2)$ . If you like, let  $x = u_1$  and  $y = u_2$ . What is  $\frac{\partial(x_1, x_2, x_3)}{\partial(x, y)}$ ?

$$D\mathbf{f}(x, y) = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 2x & 0 \end{pmatrix}$$

and so

$$D\mathbf{f}(x, y)^T = \begin{pmatrix} 1 & 0 & 2x \\ 0 & 1 & 0 \end{pmatrix}$$

Thus in this case,

$$D\mathbf{f}(\mathbf{u})^T D\mathbf{f}(\mathbf{u}) = \begin{pmatrix} 1 + 4x^2 & 0 \\ 0 & 1 \end{pmatrix}$$

and so the area element is  $\sqrt{1 + 4x^2} dx dy$  and the surface area is obtained by integrating the function,  $h(\mathbf{x}) \equiv 1$ . Therefore, this area is

$$\int_U dV = \int_0^1 \int_0^1 \sqrt{1 + 4x^2} dx dy = \frac{1}{2}\sqrt{5} - \frac{1}{4} \ln(-2 + \sqrt{5})$$

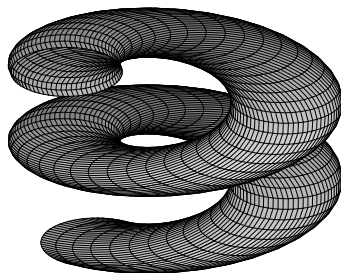
which can be obtained by using the trig. substitution,  $2x = \tan \theta$  on the inside integral.

Note this all depends on being able to write the surface in the form,  $\mathbf{x} = \mathbf{f}(\mathbf{u})$  for  $\mathbf{u} \in U \subseteq \mathbb{R}^p$ . Surfaces obtained in this form are called parametrically defined surfaces. These are best but sometimes you have some other description of a surface and in these cases things can get pretty intractable. For example, you might have a level surface of the form  $3x^2 + 4y^4 + z^6 = 10$ . In this case, you could solve for  $z$  using methods of algebra. Thus  $z = \sqrt[6]{10 - 3x^2 - 4y^4}$  and a parametric description of part of this level surface is  $(x, y, \sqrt[6]{10 - 3x^2 - 4y^4})$  for  $(x, y) \in U$  where  $U = \{(x, y) : 3x^2 + 4y^4 \leq 10\}$ . But what if the level surface was something like

$$\sin(x^2 + \ln(7 + y^2 \sin x)) + \sin(zx)e^z = 11 \sin(xyz)?$$

I really don't see how to use methods of algebra to solve for some variable in terms of the others. It isn't even clear to me whether there are any points  $(x, y, z) \in \mathbb{R}^3$  satisfying this particular relation. However, if a point satisfying this relation can be identified, the implicit function theorem from advanced calculus can usually be used to assert one of the variables is a function of the others proving the existence of a parameterization at least locally. However, this theorem doesn't give us the answer in terms of known functions so this isn't much help. Finding a parametric description of a surface is a hard problem and there are no easy answers.

**Example 22.0.38** Let  $U = [0, 12] \times [0, 2\pi]$  and let  $\mathbf{f} : U \rightarrow \mathbb{R}^3$  be given by  $\mathbf{f}(t, s) \equiv (2 \cos t + \cos s, 2 \sin t + \sin s, t)^T$ . Find a double integral for the surface area. A graph of this surface is drawn below.



It looks like something you would use to make sausages. Anyway,

$$D\mathbf{f}(t, s) = \begin{pmatrix} -2 \sin t & -\sin s \\ 2 \cos t & \cos s \\ 1 & 0 \end{pmatrix}$$

and so

$$D\mathbf{f}(t, s)^T D\mathbf{f}(t, s) = \begin{pmatrix} 5 & 2 \sin t \sin s + 2 \cos t \cos s \\ 2 \sin t \sin s + 2 \cos t \cos s & 1 \end{pmatrix}$$

and

$$\begin{aligned} \left( \frac{\partial(x_1, x_2, x_3)}{\partial(t, s)} \right)^2 &= \det \begin{pmatrix} 5 & 2 \sin t \sin s + 2 \cos t \cos s \\ 2 \sin t \sin s + 2 \cos t \cos s & 1 \end{pmatrix} \\ &= 5 - 4 \sin^2 t \sin^2 s - 8 \sin t \sin s \cos t \cos s - 4 \cos^2 t \cos^2 s \end{aligned}$$

which implies the area equals

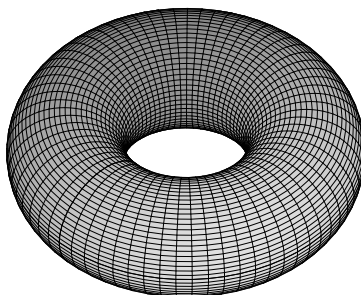
$$\int_0^{2\pi} \int_0^{12} \sqrt{5 - 4 \sin^2 t \sin^2 s - 8 \sin t \sin s \cos t \cos s - 4 \cos^2 t \cos^2 s} \, dt \, ds.$$

If you really needed to find the number this equals, how would you go about finding it? This is an interesting question and there is no single right answer. You should think about this. Here is an example for which you will be able to find the integrals.

**Example 22.0.39** Let  $U = [0, 2\pi] \times [0, 2\pi]$  and for  $(t, s) \in U$ , let

$$\mathbf{f}(t, s) = (2 \cos t + \cos t \cos s, -2 \sin t - \sin t \cos s, \sin s)^T.$$

Find the area of  $\mathbf{f}(U)$ . This is the surface of a donut shown below. The fancy name for this shape is a torus.



To find its area,

$$D\mathbf{f}(t, s) = \begin{pmatrix} -2 \sin t - \sin t \cos s & -\cos t \sin s \\ -2 \cos t - \cos t \cos s & \sin t \sin s \\ 0 & \cos s \end{pmatrix}$$

and so

$$D\mathbf{f}(t, s)^T D\mathbf{f}(t, s) = \begin{pmatrix} 4 + 4 \cos s + \cos^2 s & 0 \\ 0 & 1 \end{pmatrix}$$

which implies the area element is

$$\begin{aligned} \det \begin{pmatrix} 4 + 4 \cos s + \cos^2 s & 0 \\ 0 & 1 \end{pmatrix}^{1/2} ds dt &= (4 + 4 \cos s + \cos^2 s)^{1/2} ds dt \\ &= (\cos s + 2) ds dt \end{aligned}$$

and the area is

$$\int_0^{2\pi} \int_0^{2\pi} (\cos s + 2) ds dt = 8\pi^2$$

**Example 22.0.40** Let  $U = [0, 2\pi] \times [0, 2\pi]$  and for  $(t, s) \in U$ , let

$$\mathbf{f}(t, s) = (2 \cos t + \cos t \cos s, -2 \sin t - \sin t \cos s, \sin s)^T.$$

Find

$$\int_{\mathbf{f}(U)} h dV$$

where  $h(x, y, z) = x^2$ .

Everything is the same as the preceding example except this time it is an integral of a function. The area element is  $(\cos s + 2) ds dt$  and so the integral called for is

$$\int_{\mathbf{f}(U)} h dV = \int_0^{2\pi} \int_0^{2\pi} \left( \overbrace{2 \cos t + \cos t \cos s}^{x \text{ on the surface}} \right)^2 (\cos s + 2) ds dt = 22\pi^2$$

**Example 22.0.41** Let  $U = \{(x, y, z) : x^2 + y^2 + z^2 \leq 4\}$  and for  $(x, y, z) \in U$  let  $\mathbf{f}(x, y, z) = (x, y, x + y, z)$ . Find the three dimensional volume of  $\mathbf{f}(U)$ .



Note there is no picture here because I am unable to draw one in four dimensions. Nevertheless it is a three dimensional volume which is being computed. Everything is done the same as before.

$$D\mathbf{f}(x, y, z) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

and so

$$D\mathbf{f}(x, y, z)^T D\mathbf{f}(x, y, z) = \begin{pmatrix} 2 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

and so the volume element is  $3 \, dx \, dy \, dz$ . Therefore, the volume of  $\mathbf{f}(U)$  is

$$\int_U 3 \, dx \, dy \, dz = 3 \left( \frac{4}{3} \pi (8) \right) = 32\pi.$$

The special case where a surface is in the form  $z = f(x, y), (x, y) \in U$ , yields a simple formula which is used most often in this situation. You write the surface parametrically in the form  $\mathbf{f}(x, y) = (x, y, f(x, y))^T : (x, y) \in U$ . Then

$$D\mathbf{f}(x, y) = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ f_x & f_y \end{pmatrix}$$

and so

$$D\mathbf{f}(x, y, z)^T D\mathbf{f}(x, y, z) = \begin{pmatrix} 1 + f_x^2 & f_x f_y \\ f_x f_y & 1 + f_y^2 \end{pmatrix}.$$

Thus,

$$\det \begin{pmatrix} 1 + f_x^2 & f_x f_y \\ f_x f_y & 1 + f_y^2 \end{pmatrix} = 1 + f_y^2 + f_x^2$$

and so the area element is

$$\sqrt{1 + f_y^2 + f_x^2} \, dx \, dy.$$

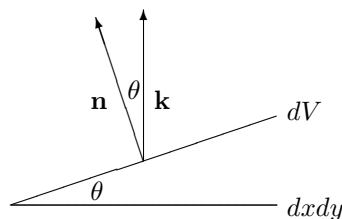
When the surface of interest comes in this simple form, people generally use this area element directly rather than worrying about a parameterization and taking determinants and finding matrices.

In the case where the surface is of the form  $x = f(y, z)$  for  $(y, z) \in U$ , the area element is obtained similarly and is

$$\sqrt{1 + f_y^2 + f_z^2} \, dy \, dz.$$

I think you can guess what the area element is if  $y = f(x, z)$ .

There is also a simple geometric description of these area elements. Consider the surface  $z = f(x, y)$ . This is a level surface of the function of three variables  $z - f(x, y)$ . In fact the surface is simply  $z - f(x, y) = 0$ . Now consider the gradient of this function of three variables. The gradient is perpendicular to the surface and the third component is positive in this case. This gradient is  $(-f_x, -f_y, 1)$  and so the unit upward normal is just  $\frac{1}{\sqrt{1 + f_x^2 + f_y^2}} (-f_x, -f_y, 1)$ . Now consider the following picture.



In this picture, you are looking at a chunk of area on the surface seen on edge and so it seems reasonable to expect to have  $dx dy = dV \cos \theta$ . But it is easy to find  $\cos \theta$  from the picture and the properties of the dot product.

$$\cos \theta = \frac{\mathbf{n} \cdot \mathbf{k}}{|\mathbf{n}| |\mathbf{k}|} = \frac{1}{\sqrt{1 + f_x^2 + f_y^2}}.$$

Therefore,  $dV = \sqrt{1 + f_x^2 + f_y^2} dx dy$  as claimed. In this context, the surface involved is referred to as  $S$  because the vector valued function,  $\mathbf{f}$  giving the parameterization will not have been identified.

**Example 22.0.42** Let  $z = \sqrt{x^2 + y^2}$  where  $(x, y) \in U$  for  $U = \{(x, y) : x^2 + y^2 \leq 4\}$  Find

$$\int_S h dV$$

where  $h(x, y, z) = x + z$  and  $S$  is the surface described as  $(x, y, \sqrt{x^2 + y^2})$  for  $(x, y) \in U$ .

Here you can see directly the angle in the above picture is  $\frac{\pi}{4}$  and so  $dV = \sqrt{2} dx dy$ . If you don't see this or if it is unclear, simply compute  $\sqrt{1 + f_x^2 + f_y^2}$  and you will find it is  $\sqrt{2}$ . Therefore, using polar coordinates,

$$\begin{aligned} \int_S h dV &= \int_U (x + \sqrt{x^2 + y^2}) \sqrt{2} dV \\ &= \sqrt{2} \int_0^{2\pi} \int_0^2 (r \cos \theta + r) r dr d\theta \\ &= \frac{16}{3} \sqrt{2} \pi. \end{aligned}$$

One other issue is worth mentioning. Suppose  $\mathbf{f}_i : U_i \rightarrow \mathbb{R}^n$  where  $U_i$  are sets in  $\mathbb{R}^p$  and suppose  $\mathbf{f}_1(U_1)$  intersects  $\mathbf{f}_2(U_2)$  along  $C$  where  $C = \mathbf{h}(V)$  for  $V \subseteq \mathbb{R}^k$  for  $k < p$ . Then define integrals and areas over  $\mathbf{f}_1(U_1) \cup \mathbf{f}_2(U_2)$  as follows.

$$\int_{\mathbf{f}_1(U_1) \cup \mathbf{f}_2(U_2)} g dV_p \equiv \int_{\mathbf{f}_1(U_1)} g dV_p + \int_{\mathbf{f}_2(U_2)} g dV_p.$$

Admittedly, the set  $C$  gets added in twice but this doesn't matter because its  $p$  dimensional volume equals zero and therefore, the integrals over this set will also be zero. Why is this? To find the  $p$  dimensional volume element on  $C$ , it is necessary to find a function,  $\mathbf{f}$ , mapping  $U \subseteq \mathbb{R}^p$  to  $C$ . Let  $\mathbf{f}(\mathbf{v}, s_1, \dots, s_{p-k}) \equiv \mathbf{h}(\mathbf{v})$ . Then  $D\mathbf{f}(\mathbf{v}, s_1, \dots, s_{p-k})$  has at least one column of zeros and so  $\det(D\mathbf{f}^T D\mathbf{f}) = 0$  showing the  $p$  dimensional volume element is zero and so this makes no contribution to the integral as claimed. Clearly something similar holds in the case of many surfaces joined in this way.

I have been purposely vague about precise mathematical conditions necessary for the above procedures. This is because the precise mathematical conditions which are usually cited are very technical and at the same time far too restrictive. The most general conditions under which these sorts of procedures are valid include things like Lipschitz functions defined on very general sets. These are functions satisfying a Lipschitz condition of the form  $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| \leq K|\mathbf{x} - \mathbf{y}|$ . For example,  $y = |x|$  is Lipschitz continuous. However, this function does not have a derivative at every point. So it is with Lipschitz functions. However, it turns out these functions have derivatives at enough points to push everything through but this requires considerations involving the Lebesgue integral. Lipschitz functions are also not the most general kind of function for which the above is valid.

## 22.1 Exercises With Answers

1. Find a parameterization for the intersection of the planes  $x + y + 2z = -3$  and  $2x - y + z = -4$ .

Answer:

$$(x, y, z) = \left(-t - \frac{7}{3}, -t - \frac{2}{3}, t\right)$$

2. Find a parameterization for the intersection of the plane  $4x + 2y + 4z = 0$  and the circular cylinder  $x^2 + y^2 = 16$ .

Answer:

The cylinder is of the form  $x = 4 \cos t, y = 4 \sin t$  and  $z = z$ . Therefore, from the equation of the plane,  $16 \cos t + 8 \sin t + 4z = 0$ . Therefore,  $z = -16 \cos t - 8 \sin t$  and this shows the parameterization is of the form  $(x, y, z) = (4 \cos t, 4 \sin t, -16 \cos t - 8 \sin t)$  where  $t \in [0, 2\pi]$ .

3. Find a parameterization for the intersection of the plane  $3x + 2y + z = 4$  and the elliptic cylinder  $x^2 + 4z^2 = 1$ .

Answer:

The cylinder is of the form  $x = \cos t, 2z = \sin t$  and  $y = y$ . Therefore, from the equation of the plane,  $3 \cos t + 2y + \frac{1}{2} \sin t = 4$ . Therefore,  $y = 2 - \frac{3}{2} \cos t - \frac{1}{4} \sin t$  and this shows the parameterization is of the form  $(x, y, z) = \left(\cos t, 2 - \frac{3}{2} \cos t - \frac{1}{4} \sin t, \frac{1}{2} \sin t\right)$  where  $t \in [0, 2\pi]$ .

4. Find a parameterization for the straight line joining  $(4, 3, 2)$  and  $(1, 7, 6)$ .

Answer:

$$(x, y, z) = (4, 3, 2) + t(-3, 4, 4) = (4 - 3t, 3 + 4t, 2 + 4t) \text{ where } t \in [0, 1].$$

5. Find a parameterization for the intersection of the surfaces  $y + 3z = 4x^2 + 4$  and  $4y + 4z = 2x + 4$ .

Answer:

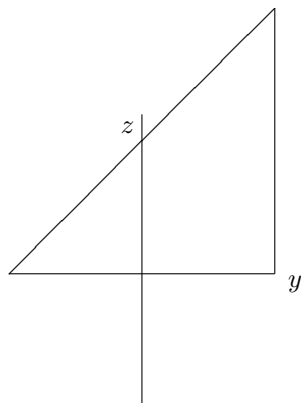
This is an application of Cramer's rule.  $y = -2x^2 - \frac{1}{2} + \frac{3}{4}x, z = -\frac{1}{4}x + \frac{3}{2} + 2x^2$ . Therefore, the parameterization is  $(x, y, z) = \left(t, -2t^2 - \frac{1}{2} + \frac{3}{4}t, -\frac{1}{4}t + \frac{3}{2} + 2t^2\right)$ .

6. Find the area of  $S$  if  $S$  is the part of the circular cylinder  $x^2 + y^2 = 16$  which lies between  $z = 0$  and  $z = 4 + y$ .

Answer:

Use the parameterization,  $x = 4 \cos v, y = 4 \sin v$  and  $z = u$  with the parameter domain described as follows. The parameter,  $v$  goes from  $-\frac{\pi}{2}$  to  $\frac{3\pi}{2}$  and for each  $v$  in

this interval,  $u$  should go from 0 to  $4 + 4 \sin v$ . To see this observe that the cylinder has its axis parallel to the  $z$  axis and if you look at a side view of the surface you would see something like this:



The positive  $x$  axis is coming out of the paper toward you in the above picture and the angle  $v$  is the usual angle measured from the positive  $x$  axis. Therefore, the area is just  $A = \int_{-\pi/2}^{3\pi/2} \int_0^{4+4\sin v} 4 \, du \, dv = 32\pi$ .

7. Find the area of  $S$  if  $S$  is the part of the cone  $x^2 + y^2 = 9z^2$  between  $z = 0$  and  $z = h$ .

Answer:

When  $z = h$ ,  $x^2 + y^2 = 9h^2$  which is the boundary of a circle of radius  $3h$ . A parameterization of this surface is  $x = u, y = v, z = \frac{1}{3}\sqrt{u^2 + v^2}$  where  $(u, v) \in D$ , a disk centered at the origin having radius  $3h$ . Therefore, the volume is just  $\int \int_D \sqrt{1 + z_u^2 + z_v^2} \, dA = \int_{-3h}^{3h} \int_{-\sqrt{9h^2 - u^2}}^{\sqrt{9h^2 - u^2}} \frac{1}{3} \sqrt{10} \, dv \, du = 3\pi h^2 \sqrt{10}$

8. Parametrizing the cylinder  $x^2 + y^2 = 4$  by  $x = 2 \cos v, y = 2 \sin v, z = u$ , show that the area element is  $dA = 2 \, du \, dv$

Answer:

It is necessary to compute  $\frac{\partial(x,y,z)}{\partial(u,v)} = \det(D\mathbf{f}^T D\mathbf{f})$ .

$$D\mathbf{f}(u, v) = \begin{pmatrix} 0 & -2 \sin v \\ 0 & 2 \cos v \\ 1 & 0 \end{pmatrix}$$

and so  $D\mathbf{f}^T D\mathbf{f} = \begin{pmatrix} 1 & 0 \\ 0 & 4 \end{pmatrix}$  and so the area element is as described.

9. Find the area enclosed by the limaçon  $r = 2 + \cos \theta$ .

Answer:

You can graph this region and you see it is sort of an oval shape and that  $\theta \in [0, 2\pi]$  while  $r$  goes from 0 up to  $2 + \cos \theta$ . Now  $x = r \cos \theta$  and  $y = r \sin \theta$  are the  $x$  and  $y$  coordinates corresponding to  $r$  and  $\theta$  in the above parameter domain. Therefore, the area of the limaçon equals  $\int \int_P \left| \frac{\partial(x,y)}{\partial(r,\theta)} \right| \, dr \, d\theta = \int_0^{2\pi} \int_0^{2+\cos \theta} r \, dr \, d\theta$  because the Jacobian equals  $r$  in this case. Therefore, the area equals  $\int_0^{2\pi} \int_0^{2+\cos \theta} r \, dr \, d\theta = \frac{9}{2}\pi$ .

10. Find the surface area of the paraboloid  $z = h(1 - x^2 - y^2)$  between  $z = 0$  and  $z = h$ .

Answer:

Let  $R$  denote the unit circle. Then the area of the surface above this circle would be  $\int \int_R \sqrt{1 + 4x^2h^2 + 4y^2h^2} dA$ . Changing to polar coordinates, this becomes

$$\int_0^{2\pi} \int_0^1 \sqrt{1 + 4h^2r^2} r dr d\theta = \frac{1}{6}\pi \frac{\sqrt{(1+4h^2)+4}\sqrt{(1+4h^2)h^2-1}}{h^2}.$$

11. Evaluate  $\int \int_S (1+x) dA$  where  $S$  is the part of the plane  $2x + 3y + 3z = 18$  which is in the first octant.

Answer:

$$\int_0^6 \int_0^{6-\frac{2}{3}x} (1+x) \frac{1}{3}\sqrt{22} dy dx = 28\sqrt{22}$$

12. Evaluate  $\int \int_S (1+x) dA$  where  $S$  is the part of the cylinder  $x^2 + y^2 = 16$  between  $z = 0$  and  $z = h$ .

Answer:

Parametrize the cylinder as  $x = 4\cos\theta$  and  $y = 4\sin\theta$  while  $z = t$  and the parameter domain is just  $[0, 2\pi] \times [0, h]$ . Then the integral to evaluate would be

$$\int_0^{2\pi} \int_0^h (1 + 4\cos\theta) 4 dt d\theta = 8h\pi.$$

Note how  $4\cos\theta$  was substituted for  $x$  and the area element is  $4 dt d\theta$ .

13. Evaluate  $\int \int_S (1+x) dA$  where  $S$  is the hemisphere  $x^2 + y^2 + z^2 = 16$  between  $x = 0$  and  $x = 4$ .

Answer:

Parametrize the sphere as  $x = 4\sin\phi\cos\theta$ ,  $y = 4\sin\phi\sin\theta$ , and  $z = 4\cos\phi$  and consider the values of the parameters. Since it is referred to as a hemisphere and involves  $x > 0$ ,  $\theta \in [-\frac{\pi}{2}, \frac{\pi}{2}]$  and  $\phi \in [0, \pi]$ . Then the area element is  $\sqrt{a^4\sin\phi} d\theta d\phi$  and so the integral to evaluate is

$$\int_0^\pi \int_{-\pi/2}^{\pi/2} (1 + 4\sin\phi\cos\theta) 16\sin\phi d\theta d\phi = 96\pi$$

14. For  $(\theta, \alpha) \in [0, 2\pi] \times [0, 2\pi]$ , let

$$\mathbf{f}(\theta, \alpha) \equiv (\cos\theta(2 + \cos\alpha), -\sin\theta(2 + \cos\alpha), \sin\alpha)^T.$$

Find the area of  $\mathbf{f}([0, 2\pi] \times [0, 2\pi])$ .

Answer:

$$D\mathbf{f}(\theta, \alpha) = \begin{pmatrix} -\sin(\theta)(2 + \cos\alpha) & -\cos\theta\sin\alpha \\ -\cos(\theta)(2 + \cos\alpha) & \sin\theta\sin\alpha \\ 0 & \cos\alpha \end{pmatrix} \text{ and so the area element is}$$

$$\det(D\mathbf{f}^T D\mathbf{f})^{1/2} d\theta d\alpha = (4 + 4\cos\alpha + \cos^2\alpha)^{1/2} d\theta d\alpha.$$

Therefore, the area is

$$\int_0^{2\pi} \int_0^{2\pi} (4 + 4\cos\alpha + \cos^2\alpha)^{1/2} d\theta d\alpha = \int_0^{2\pi} \int_0^{2\pi} (2 + \cos\alpha) d\theta d\alpha = 8\pi^2.$$

15. For  $(\theta, \alpha) \in [0, 2\pi] \times [0, 2\pi]$ , let

$$\mathbf{f}(\theta, \alpha) \equiv (\cos \theta (4 + 2 \cos \alpha), -\sin \theta (4 + 2 \cos \alpha), 2 \sin \alpha)^T.$$

Also let  $h(\mathbf{x}) = \cos \alpha$  where  $\alpha$  is such that

$$\mathbf{x} = (\cos \theta (4 + 2 \cos \alpha), -\sin \theta (4 + 2 \cos \alpha), 2 \sin \alpha)^T.$$

Find  $\int_{\mathbf{f}([0, 2\pi] \times [0, 2\pi])} h \, dV$ .

Answer:

$$D\mathbf{f}(\theta, \alpha) = \begin{pmatrix} -\sin(\theta)(4 + 2 \cos \alpha) & -2 \cos \theta \sin \alpha \\ -\cos(\theta)(4 + 2 \cos \alpha) & 2 \sin \theta \sin \alpha \\ 0 & 2 \cos \alpha \end{pmatrix} \text{ and so the area element is}$$

$$\det(D\mathbf{f}^T D\mathbf{f})^{1/2} \, d\theta \, d\alpha = (64 + 64 \cos \alpha + 16 \cos^2 \alpha)^{1/2} \, d\theta \, d\alpha.$$

Therefore, the desired integral is

$$\begin{aligned} & \int_0^{2\pi} \int_0^{2\pi} (\cos \alpha) (64 + 64 \cos \alpha + 16 \cos^2 \alpha)^{1/2} \, d\theta \, d\alpha \\ &= \int_0^{2\pi} \int_0^{2\pi} (\cos \alpha) (8 + 4 \cos \alpha) \, d\theta \, d\alpha = 8\pi^2 \end{aligned}$$

16. For  $(\theta, \alpha) \in [0, 2\pi] \times [0, 2\pi]$ , let

$$\mathbf{f}(\theta, \alpha) \equiv (\cos \theta (3 + \cos \alpha), -\sin \theta (3 + \cos \alpha), \sin \alpha)^T.$$

Also let  $h(\mathbf{x}) = \cos^2 \theta$  where  $\theta$  is such that

$$\mathbf{x} = (\cos \theta (3 + \cos \alpha), -\sin \theta (3 + \cos \alpha), \sin \alpha)^T.$$

Find  $\int_{\mathbf{f}([0, 2\pi] \times [0, 2\pi])} h \, dV$ .

Answer:

$$D\mathbf{f}(\theta, \alpha) = \begin{pmatrix} -\sin(\theta)(3 + \cos \alpha) & -\cos \theta \sin \alpha \\ -\cos(\theta)(3 + \cos \alpha) & \sin \theta \sin \alpha \\ 0 & \cos \alpha \end{pmatrix} \text{ and so the area element is}$$

$$\det(D\mathbf{f}^T D\mathbf{f})^{1/2} \, d\theta \, d\alpha = (9 + 6 \cos \alpha + \cos^2 \alpha)^{1/2} \, d\theta \, d\alpha.$$

Therefore, the desired integral is

$$\begin{aligned} & \int_0^{2\pi} \int_0^{2\pi} (\cos^2 \theta) (9 + 6 \cos \alpha + \cos^2 \alpha)^{1/2} \, d\theta \, d\alpha \\ &= \int_0^{2\pi} \int_0^{2\pi} (\cos^2 \theta) (3 + \cos \alpha) \, d\theta \, d\alpha = 6\pi^2 \end{aligned}$$

17. For  $(\theta, \alpha) \in [0, 25] \times [0, 2\pi]$ , let

$$\mathbf{f}(\theta, \alpha) \equiv (\cos \theta (4 + 2 \cos \alpha), -\sin \theta (4 + 2 \cos \alpha), 2 \sin \alpha + \theta)^T.$$

Find a double integral which gives the area of  $\mathbf{f}([0, 25] \times [0, 2\pi])$ .

Answer:

In this case,  $D\mathbf{f}(\theta, \alpha) = \begin{pmatrix} -\sin(\theta)(4 + 2\cos\alpha) & -2\cos\theta\sin\alpha \\ -\cos(\theta)(4 + 2\cos\alpha) & 2\sin\theta\sin\alpha \\ 1 & 2\cos\alpha \end{pmatrix}$  and so the area

element is

$\det(D\mathbf{f}^T D\mathbf{f}) d\theta d\alpha = (68 + 64\cos\alpha + 12\cos^2\alpha)^{1/2} d\theta d\alpha$  and so the surface area is

$$\int_0^{2\pi} \int_0^{2\pi} (68 + 64\cos\alpha + 12\cos^2\alpha)^{1/2} d\theta d\alpha.$$

18. For  $(\theta, \alpha) \in [0, 2\pi] \times [0, 2\pi]$ , and  $\beta$  a fixed real number, define  $\mathbf{f}(\theta, \alpha) \equiv$

$$(\cos\theta(2 + \cos\alpha), -\cos\beta\sin\theta(2 + \cos\alpha) + \sin\beta\sin\alpha, \sin\beta\sin\theta(2 + \cos\alpha) + \cos\beta\sin\alpha)^T.$$

Find a double integral which gives the area of  $\mathbf{f}([0, 2\pi] \times [0, 2\pi])$ .

Answer:

$D\mathbf{f} =$

$$\begin{pmatrix} -\sin(2\theta + \theta\cos\alpha) & -\sin\alpha\cos\theta \\ -2\cos\beta\cos\theta - \cos\beta\cos\theta\cos\alpha & \cos\beta\sin\theta\sin\alpha + \sin\beta\cos\alpha \\ 2\sin\beta\cos\theta + \sin\beta\cos\theta\cos\alpha & -\sin(\beta\sin\theta\sin\alpha) + \cos\beta\cos\alpha \end{pmatrix}$$
 and so after many

computations, the area element is  $(4 + 4\cos\alpha + \cos^2\alpha)^{1/2} d\theta d\alpha$ . Therefore, the area is  $\int_0^{2\pi} \int_0^{2\pi} (2 + \cos\alpha) d\theta d\alpha = 8\pi^2$ .

## 22.2 Exercises

- Find a parameterization for the intersection of the planes  $4x + 2y + 4z = 3$  and  $6x - 2y = -1$ .
- Find a parameterization for the intersection of the plane  $3x + y + z = 1$  and the circular cylinder  $x^2 + y^2 = 1$ .
- Find a parameterization for the intersection of the plane  $3x + 2y + 4z = 4$  and the elliptic cylinder  $x^2 + 4z^2 = 16$ .
- Find a parameterization for the straight line joining  $(1, 3, 1)$  and  $(-2, 5, 3)$ .
- Find a parameterization for the intersection of the surfaces  $4y + 3z = 3x^2 + 2$  and  $3y + 2z = -x + 3$ .
- Find the area of  $S$  if  $S$  is the part of the circular cylinder  $x^2 + y^2 = 4$  which lies between  $z = 0$  and  $z = 2 + y$ .
- Find the area of  $S$  if  $S$  is the part of the cone  $x^2 + y^2 = 16z^2$  between  $z = 0$  and  $z = h$ .
- Parametrizing the cylinder  $x^2 + y^2 = a^2$  by  $x = a\cos v, y = a\sin v, z = u$ , show that the area element is  $dA = a du dv$ .
- Find the area enclosed by the limaçon  $r = 2 + \cos\theta$ .
- Find the surface area of the paraboloid  $z = h(1 - x^2 - y^2)$  between  $z = 0$  and  $z = h$ .
- Evaluate  $\int \int_S (1 + x) dA$  where  $S$  is the part of the plane  $4x + y + 3z = 12$  which is in the first octant.

12. Evaluate  $\int_S (1+x) dA$  where  $S$  is the part of the cylinder  $x^2 + y^2 = 9$  between  $z = 0$  and  $z = h$ .
13. Evaluate  $\int_S (1+x) dA$  where  $S$  is the hemisphere  $x^2 + y^2 + z^2 = 4$  between  $x = 0$  and  $x = 2$ .
14. For  $(\theta, \alpha) \in [0, 2\pi] \times [0, 2\pi]$ , let  $\mathbf{f}(\theta, \alpha) \equiv (\cos \theta (4 + \cos \alpha), -\sin \theta (4 + \cos \alpha), \sin \alpha)^T$ . Find the area of  $\mathbf{f}([0, 2\pi] \times [0, 2\pi])$ .
15. For  $(\theta, \alpha) \in [0, 2\pi] \times [0, 2\pi]$ , let  $\mathbf{f}(\theta, \alpha) \equiv$

$$(\cos \theta (3 + 2 \cos \alpha), -\sin \theta (3 + 2 \cos \alpha), 2 \sin \alpha)^T.$$

Also let  $h(\mathbf{x}) = \cos \alpha$  where  $\alpha$  is such that

$$\mathbf{x} = (\cos \theta (3 + 2 \cos \alpha), -\sin \theta (3 + 2 \cos \alpha), 2 \sin \alpha)^T.$$

Find  $\int_{\mathbf{f}([0, 2\pi] \times [0, 2\pi])} h dV$ .

16. For  $(\theta, \alpha) \in [0, 2\pi] \times [0, 2\pi]$ , let  $\mathbf{f}(\theta, \alpha) \equiv$

$$(\cos \theta (4 + 3 \cos \alpha), -\sin \theta (4 + 3 \cos \alpha), 3 \sin \alpha)^T.$$

Also let  $h(\mathbf{x}) = \cos^2 \theta$  where  $\theta$  is such that

$$\mathbf{x} = (\cos \theta (4 + 3 \cos \alpha), -\sin \theta (4 + 3 \cos \alpha), 3 \sin \alpha)^T.$$

Find  $\int_{\mathbf{f}([0, 2\pi] \times [0, 2\pi])} h dV$ .

17. For  $(\theta, \alpha) \in [0, 2\pi] \times [0, 2\pi]$ , let  $\mathbf{f}(\theta, \alpha) \equiv$

$$(\cos \theta (4 + 2 \cos \alpha), -\sin \theta (4 + 2 \cos \alpha), 2 \sin \alpha + \theta)^T.$$

Find a double integral which gives the area of  $\mathbf{f}([0, 2\pi] \times [0, 2\pi])$ .

18. For  $(\theta, \alpha) \in [0, 2\pi] \times [0, 2\pi]$ , and  $\beta$  a fixed real number, define  $\mathbf{f}(\theta, \alpha) \equiv$

$$\begin{pmatrix} \cos \theta (3 + 2 \cos \alpha), -\cos \beta \sin \theta (3 + 2 \cos \alpha) + \\ 2 \sin \beta \sin \alpha, \sin \beta \sin \theta (3 + 2 \cos \alpha) + 2 \cos \beta \sin \alpha \end{pmatrix}^T.$$

Find a double integral which gives the area of  $\mathbf{f}([0, 2\pi] \times [0, 2\pi])$ .

19. In the case where  $\mathbf{f} : U \rightarrow \mathbb{R}^3$ , show that

$$\frac{\partial (x, y, z)}{\partial (u_1, u_2)} = |\mathbf{f}_{u_1} \times \mathbf{f}_{u_2}|.$$

Thus the area element is  $|\mathbf{f}_{u_1} \times \mathbf{f}_{u_2}| du_1 du_2$ .

20. In spherical coordinates,  $\phi = c, \rho \in [0, R]$  determines a cone. Find the area of this cone without doing any work involving Jacobians and such.



# Vector Calculus

## 23.1 Divergence And Curl Of A Vector Field

Imagine a region,  $U$  in  $\mathbb{R}^3$  and at each point,  $\mathbf{x}$ , of  $U$  there is associated a vector,  $\mathbf{f}(\mathbf{x})$ . For example, this could be the force acting on a unit mass at the point,  $\mathbf{x}$ . It could be the velocity of a fluid moving through  $\mathbb{R}^3$  at  $\mathbf{x}$ , etc. Such a thing is called a vector field. The precise definition follows.

**Definition 23.1.1** A function,  $\mathbf{f} : U \rightarrow \mathbb{R}^p$  for  $U \subseteq \mathbb{R}^p$  an open set is called a vector field. A scalar valued function is called a scalar field. It is called a  $C^k$  vector field if the function,  $\mathbf{f}$  is a  $C^k$  function. For a  $C^1$  vector field, as just described  $\nabla \cdot \mathbf{f}(\mathbf{x}) \equiv \text{div } \mathbf{f}(\mathbf{x})$  known as the divergence, is defined as

$$\nabla \cdot \mathbf{f}(\mathbf{x}) \equiv \text{div } \mathbf{f}(\mathbf{x}) \equiv \sum_{i=1}^p \frac{\partial f_i}{\partial x_i}(\mathbf{x}).$$

Using the repeated summation convention, this is often written as

$$f_{i,i}(\mathbf{x}) \equiv \partial_i f_i(\mathbf{x})$$

where the comma indicates a partial derivative is being taken with respect to the  $i^{\text{th}}$  variable and  $\partial_i$  denotes differentiation with respect to the  $i^{\text{th}}$  variable. In words, the divergence is the sum of the  $i^{\text{th}}$  derivative of the  $i^{\text{th}}$  component function of  $\mathbf{f}$  for all values of  $i$ . If  $p = 3$ , the curl of the vector field yields another vector field and it is defined as follows.

$$(\text{curl } (\mathbf{f})(\mathbf{x}))_i \equiv (\nabla \times \mathbf{f}(\mathbf{x}))_i \equiv \varepsilon_{ijk} \partial_j f_k(\mathbf{x})$$

where here  $\partial_j$  means the partial derivative with respect to  $x_j$  and the subscript of  $i$  in  $(\text{curl } (\mathbf{f})(\mathbf{x}))_i$  means the  $i^{\text{th}}$  Cartesian component of the vector,  $\text{curl } (\mathbf{f})(\mathbf{x})$ . Thus the curl is evaluated by expanding the following determinant along the top row.

$$\begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ \frac{\partial}{\partial x} & \frac{\partial}{\partial y} & \frac{\partial}{\partial z} \\ f_1(x, y, z) & f_2(x, y, z) & f_3(x, y, z) \end{vmatrix}.$$

Note the similarity with the cross product. Sometimes the curl is called rot. (Short for rotation not decay.) Also

$$\nabla^2 f \equiv \nabla \cdot (\nabla f).$$

This last symbol is important enough that it is given a name, the Laplacian. It is also denoted by  $\Delta$ . Thus  $\nabla^2 f = \Delta f$ . In addition for  $\mathbf{f}$  a vector field, the symbol  $\mathbf{f} \cdot \nabla$  is defined as a

“differential operator” in the following way.

$$\mathbf{f} \cdot \nabla (\mathbf{g}) \equiv f_1(\mathbf{x}) \frac{\partial \mathbf{g}(\mathbf{x})}{\partial x_1} + f_2(\mathbf{x}) \frac{\partial \mathbf{g}(\mathbf{x})}{\partial x_2} + \cdots + f_p(\mathbf{x}) \frac{\partial \mathbf{g}(\mathbf{x})}{\partial x_p}.$$

Thus  $\mathbf{f} \cdot \nabla$  takes vector fields and makes them into new vector fields.

This definition is in terms of a given coordinate system but later coordinate free definitions of the curl and div are presented. For now, everything is defined in terms of a given Cartesian coordinate system. The divergence and curl have profound physical significance and this will be discussed later. For now it is important to understand their definition in terms of coordinates. Be sure you understand that for  $\mathbf{f}$  a vector field,  $\text{div } \mathbf{f}$  is a scalar field meaning it is a scalar valued function of three variables. For a scalar field,  $f$ ,  $\nabla f$  is a vector field described earlier on Page 493. For  $\mathbf{f}$  a vector field having values in  $\mathbb{R}^3$ ,  $\text{curl } \mathbf{f}$  is another vector field.

**Example 23.1.2** Let  $\mathbf{f}(\mathbf{x}) = xy\mathbf{i} + (z - y)\mathbf{j} + (\sin(x) + z)\mathbf{k}$ . Find  $\text{div } \mathbf{f}$  and  $\text{curl } \mathbf{f}$ .

First the divergence of  $\mathbf{f}$  is

$$\frac{\partial(xy)}{\partial x} + \frac{\partial(z - y)}{\partial y} + \frac{\partial(\sin(x) + z)}{\partial z} = y + (-1) + 1 = y.$$

Now  $\text{curl } \mathbf{f}$  is obtained by evaluating

$$\begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ \frac{\partial}{\partial x} & \frac{\partial}{\partial y} & \frac{\partial}{\partial z} \\ xy & z - y & \sin(x) + z \end{vmatrix} =$$

$$\mathbf{i} \left( \frac{\partial}{\partial y} (\sin(x) + z) - \frac{\partial}{\partial z} (z - y) \right) - \mathbf{j} \left( \frac{\partial}{\partial x} (\sin(x) + z) - \frac{\partial}{\partial z} (xy) \right) +$$

$$\mathbf{k} \left( \frac{\partial}{\partial x} (z - y) - \frac{\partial}{\partial y} (xy) \right) = -\mathbf{i} - \cos(x)\mathbf{j} - x\mathbf{k}.$$

There are many interesting identities which relate the gradient, divergence and curl.

**Theorem 23.1.3** Assuming  $\mathbf{f}, \mathbf{g}$  are a  $C^2$  vector fields whenever necessary, the following identities are valid.

1.  $\nabla \cdot (\nabla \times \mathbf{f}) = 0$
2.  $\nabla \times \nabla \phi = \mathbf{0}$
3.  $\nabla \times (\nabla \times \mathbf{f}) = \nabla(\nabla \cdot \mathbf{f}) - \nabla^2 \mathbf{f}$  where  $\nabla^2 \mathbf{f}$  is a vector field whose  $i^{\text{th}}$  component is  $\nabla^2 f_i$ .
4.  $\nabla \cdot (\mathbf{f} \times \mathbf{g}) = \mathbf{g} \cdot (\nabla \times \mathbf{f}) - \mathbf{f} \cdot (\nabla \times \mathbf{g})$
5.  $\nabla \times (\mathbf{f} \times \mathbf{g}) = (\nabla \cdot \mathbf{g})\mathbf{f} - (\nabla \cdot \mathbf{f})\mathbf{g} + (\mathbf{g} \cdot \nabla)\mathbf{f} - (\mathbf{f} \cdot \nabla)\mathbf{g}$

**Proof:** These are all easy to establish if you use the repeated index summation convention and the reduction identities discussed on Page 331.

$$\begin{aligned}
 \nabla \cdot (\nabla \times \mathbf{f}) &= \partial_i (\nabla \times \mathbf{f})_i \\
 &= \partial_i (\varepsilon_{ijk} \partial_j f_k) \\
 &= \varepsilon_{ijk} \partial_i (\partial_j f_k) \\
 &= \varepsilon_{jik} \partial_j (\partial_i f_k) \\
 &= -\varepsilon_{ijk} \partial_j (\partial_i f_k) \\
 &= -\varepsilon_{ijk} \partial_i (\partial_j f_k) \\
 &= -\nabla \cdot (\nabla \times \mathbf{f}).
 \end{aligned}$$

This establishes the first formula. The second formula is done similarly. Now consider the third.

$$\begin{aligned}
 (\nabla \times (\nabla \times \mathbf{f}))_i &= \varepsilon_{ijk} \partial_j (\nabla \times \mathbf{f})_k \\
 &= \varepsilon_{ijk} \partial_j (\varepsilon_{krs} \partial_r f_s) \\
 &= \varepsilon_{ijk} \overbrace{\varepsilon_{krs}}^{\varepsilon_{kij}} \partial_j (\partial_r f_s) \\
 &= (\delta_{ir} \delta_{js} - \delta_{is} \delta_{jr}) \partial_j (\partial_r f_s) \\
 &= \partial_j (\partial_i f_j) - \partial_j (\partial_j f_i) \\
 &= \partial_i (\partial_j f_j) - \partial_j (\partial_j f_i) \\
 &= (\nabla (\nabla \cdot \mathbf{f}) - \nabla^2 \mathbf{f})_i
 \end{aligned}$$

This establishes the third identity.

Consider the fourth identity.

$$\begin{aligned}
 \nabla \cdot (\mathbf{f} \times \mathbf{g}) &= \partial_i (\mathbf{f} \times \mathbf{g})_i \\
 &= \partial_i \varepsilon_{ijk} f_j g_k \\
 &= \varepsilon_{ijk} (\partial_i f_j) g_k + \varepsilon_{ijk} f_j (\partial_i g_k) \\
 &= (\varepsilon_{kij} \partial_i f_j) g_k - (\varepsilon_{jik} \partial_i g_k) f_k \\
 &= \nabla \times \mathbf{f} \cdot \mathbf{g} - \nabla \times \mathbf{g} \cdot \mathbf{f}.
 \end{aligned}$$

This proves the fourth identity.

Consider the fifth.

$$\begin{aligned}
 (\nabla \times (\mathbf{f} \times \mathbf{g}))_i &= \varepsilon_{ijk} \partial_j (\mathbf{f} \times \mathbf{g})_k \\
 &= \varepsilon_{ijk} \partial_j \varepsilon_{krs} f_r g_s \\
 &= \varepsilon_{kij} \varepsilon_{krs} \partial_j (f_r g_s) \\
 &= (\delta_{ir} \delta_{js} - \delta_{is} \delta_{jr}) \partial_j (f_r g_s) \\
 &= \partial_j (f_i g_j) - \partial_j (f_j g_i) \\
 &= (\partial_j g_j) f_i + g_j \partial_j f_i - (\partial_j f_j) g_i - f_j (\partial_j g_i) \\
 &= ((\nabla \cdot \mathbf{g}) \mathbf{f} + (\mathbf{g} \cdot \nabla) \mathbf{f}) - ((\nabla \cdot \mathbf{f}) \mathbf{g} + (\mathbf{f} \cdot \nabla) \mathbf{g})_i
 \end{aligned}$$

and this establishes the fifth identity.

I think the important thing about the above is not that these identities can be proved and are valid as much as the method by which they were proved. The reduction identities on Page 331 were used to discover the identities. There is a difference between proving something someone tells you about and both discovering what should be proved and proving it. This

notation and the reduction identity make the discovery of vector identities fairly routine and this is why these things are of great significance.

One of the above identities says  $\nabla \cdot (\nabla \times \mathbf{f}) = 0$ . Suppose now  $\nabla \cdot \mathbf{g} = 0$ . Does it follow that there exists  $\mathbf{f}$  such that  $\mathbf{g} = \nabla \times \mathbf{f}$ ? It turns out that this is usually the case and when such an  $\mathbf{f}$  exists, it is called a vector potential. Here is one way to do it, assuming everything is defined so the following formulas make sense.

$$\mathbf{f}(x, y, z) = \left( \int_0^z g_2(x, y, t) dt, - \int_0^z g_1(x, y, t) dt + \int_0^x g_3(t, y, 0) dt, 0 \right)^T. \quad (23.1)$$

In verifying this you need to use the following manipulation which will generally hold under reasonable conditions but which has not been carefully shown yet.

$$\frac{\partial}{\partial x} \int_a^b h(x, t) dt = \int_a^b \frac{\partial h}{\partial x}(x, t) dt. \quad (23.2)$$

The above formula seems plausible because the integral is a sort of a sum and the derivative of a sum is the sum of the derivatives. However, this sort of sloppy reasoning will get you into all sorts of trouble. The formula involves the interchange of two limit operations, the integral and the limit of a difference quotient. Such an interchange can only be accomplished through a theorem. The following gives the necessary result.

**Lemma 23.1.4** Suppose  $h$  and  $\frac{\partial h}{\partial x}$  are continuous on the rectangle  $R = [c, d] \times [a, b]$ . Then (23.2) holds.

**Proof:** Let  $H(x) = \int_a^b \frac{\partial h}{\partial x}(x, t) dt$ .

**Claim:**  $H$  is continuous.

From Theorem 14.9.8 on Page 348,  $\frac{\partial h}{\partial x}$  is uniformly continuous on  $R$  since  $R$  is a closed and bounded set. Therefore, letting  $\varepsilon > 0$  be given, there exists  $\delta > 0$  such that if  $|x - x'| < \delta$ , then

$$\left| \frac{\partial h}{\partial x}(x, t) - \frac{\partial h}{\partial x}(x', t) \right| < \frac{\varepsilon}{b - a + 1}.$$

for all  $t \in [a, b]$ . Consequently,

$$\begin{aligned} |H(x) - H(x')| &= \left| \int_a^b \frac{\partial h}{\partial x}(x, t) dt - \int_a^b \frac{\partial h}{\partial x}(x', t) dt \right| \\ &\leq \int_a^b \left| \frac{\partial h}{\partial x}(x, t) - \frac{\partial h}{\partial x}(x', t) \right| dt < \frac{\varepsilon}{b - a + 1} (b - a) < \varepsilon. \end{aligned}$$

Since  $\varepsilon > 0$  is arbitrary, this establishes the claim.

Now by Fubini's theorem,

$$\begin{aligned} \int_c^x H(r) dr &= \int_c^x \int_a^b \frac{\partial h}{\partial x}(r, t) dt dr \\ &= \int_a^b \int_c^x \frac{\partial h}{\partial x}(r, t) dr dt \\ &= \int_a^b h(x, t) dt - \int_a^b h(c, t) dt. \end{aligned}$$

Using the fundamental theorem of calculus in the above formula, it follows

$$\frac{\partial}{\partial x} \left( \int_a^b h(x, t) dt - \int_a^b h(c, t) dt \right)$$

exists and equals  $H(x)$ . Thus

$$\begin{aligned} H(x) &\equiv \int_a^b \frac{\partial h}{\partial x}(x, t) dt = \frac{\partial}{\partial x} \left( \int_a^b h(x, t) dt - \int_a^b h(c, t) dt \right) \\ &= \frac{\partial}{\partial x} \left( \int_a^b h(x, t) dt \right) \end{aligned}$$

and this proves the lemma.

The second formula of Theorem 23.1.3 states  $\nabla \times \nabla \phi = \mathbf{0}$ . This suggests the following question: Suppose  $\nabla \times \mathbf{f} = \mathbf{0}$ , does it follow there exists  $\phi$ , a scalar field such that  $\nabla \phi = \mathbf{f}$ ? The answer to this is often yes and a theorem will be given and proved after the presentation of Stoke's theorem. This scalar field,  $\phi$ , is called a scalar potential for  $\mathbf{f}$ .

There is also a fundamental result having great significance which involves  $\nabla^2$  called the maximum principle. This principle says that if  $\nabla^2 u \geq 0$  on a bounded open set,  $U$ , then  $u$  achieves its maximum value on the boundary of  $U$ . You must understand the fundamental advanced calculus results involving compactness in order to read the proof of this theorem.

**Theorem 23.1.5** *Let  $U$  be a bounded open set in  $\mathbb{R}^n$  and suppose  $u \in C^2(U) \cap C(\overline{U})$  such that  $\nabla^2 u \geq 0$  in  $U$ . Then letting  $\partial U = \overline{U} \setminus U$ , it follows that  $\max \{u(\mathbf{x}) : \mathbf{x} \in \overline{U}\} = \max \{u(\mathbf{x}) : \mathbf{x} \in \partial U\}$ .*

**Proof:** If this is not so, there exists  $\mathbf{x}_0 \in U$  such that  $u(\mathbf{x}_0) > \max \{u(\mathbf{x}) : \mathbf{x} \in \partial U\}$ . Since  $U$  is bounded there exists  $\varepsilon > 0$  such that  $u(\mathbf{x}) + \varepsilon |\mathbf{x}|^2$  also has its maximum in  $U$  because if this is not so, there exists a sequence,  $\{\varepsilon_n\}$  of positive numbers converging to zero and a point  $\mathbf{x}_{\varepsilon_n} \in \partial U$  such that  $u(\mathbf{x}_{\varepsilon_n}) + \varepsilon_n |\mathbf{x}_{\varepsilon_n}|^2 \geq u(\mathbf{x}) + \varepsilon_n |\mathbf{x}|^2$  for all  $\mathbf{x} \in \overline{U}$ . Then using compactness of  $\partial U$ , there exists a subsequence, still denoted by  $\varepsilon_n$  such that  $\mathbf{x}_{\varepsilon_n} \rightarrow \mathbf{x}_1 \in \partial U$  and so, taking the limit,

$$u(\mathbf{x}_1) \geq u(\mathbf{x}) \text{ for all } \mathbf{x} \in \overline{U},$$

contrary to what was assumed about  $\mathbf{x}_0$ .

Now let  $\mathbf{x}_1$  be the point in  $U$  at which  $u(\mathbf{x}) + \varepsilon |\mathbf{x}|^2$  achieves its maximum. As an exercise you should show that  $\nabla^2(f + g) = \nabla^2 f + \nabla^2 g$  and therefore,  $\nabla^2(u(\mathbf{x}) + \varepsilon |\mathbf{x}|^2) = \nabla^2 u(\mathbf{x}) + 2n\varepsilon$ . (Why?) Therefore,

$$0 \geq \nabla^2 u(\mathbf{x}_1) + 2n\varepsilon \geq 2n\varepsilon,$$

a contradiction. This proves the theorem.

## 23.2 Exercises

1. Find  $\operatorname{div} \mathbf{f}$  and  $\operatorname{curl} \mathbf{f}$  where  $\mathbf{f}$  is

- (a)  $(xyz, x^2 + \ln(xy), \sin x^2 + z)^T$
- (b)  $(\sin x, \sin y, \sin z)^T$
- (c)  $(f(x), g(y), h(z))^T$
- (d)  $(x - 2, y - 3, z - 6)^T$
- (e)  $(y^2, 2xy, \cos z)^T$

$$(f) \ (f(y, z), g(x, z), h(y, z))^T$$

2. Prove formula 2 of Theorem 23.1.3.
3. Simplify the expression  $\mathbf{f} \times (\nabla \times \mathbf{g}) + \mathbf{g} \times (\nabla \times \mathbf{f}) + (\mathbf{f} \cdot \nabla) \mathbf{g} + (\mathbf{g} \cdot \nabla) \mathbf{f}$ .
4. Simplify  $\nabla \times (\mathbf{v} \times \mathbf{r})$  where  $\mathbf{r} = (x, y, z)^T = x\mathbf{i} + y\mathbf{j} + z\mathbf{k}$  and  $\mathbf{v}$  is a constant vector.
5. Discover a formula which simplifies  $\nabla \cdot (v \nabla u)$ .
6. Verify that  $\nabla \cdot (u \nabla v) - \nabla \cdot (v \nabla u) = u \nabla^2 v - v \nabla^2 u$ .
7. Verify that  $\nabla^2 (uv) = v \nabla^2 u + 2(\nabla u \cdot \nabla v) + u \nabla^2 v$ .
8. Functions,  $u$ , which satisfy  $\nabla^2 u = 0$  are called harmonic functions. Show the following functions are harmonic wherever they are defined.

- (a)  $2xy$
- (b)  $x^2 - y^2$
- (c)  $\sin x \cosh y$
- (d)  $\ln(x^2 + y^2)$
- (e)  $1/\sqrt{x^2 + y^2 + z^2}$

9. Verify the formula given in (23.1) is a vector potential for  $\mathbf{g}$  assuming that  $\operatorname{div} \mathbf{g} = 0$ .
10. Show that if  $\nabla^2 u_k = 0$  for each  $k = 1, 2, \dots, m$ , and  $c_k$  is a constant, then  $\nabla^2 (\sum_{k=1}^m c_k u_k) = 0$  also.
11. In Theorem 23.1.5 why is  $\nabla^2 (\varepsilon |\mathbf{x}|^2) = 2n\varepsilon$ ?
12. Using Theorem 23.1.5 prove the following: Let  $f \in C(\partial U)$  ( $f$  is continuous on  $\partial U$ ) where  $U$  is a bounded open set. Then there exists at most one solution,  $u \in C^2(U) \cap C(\overline{U})$  and  $\nabla^2 u = 0$  in  $U$  with  $u = f$  on  $\partial U$ . **Hint:** Suppose there are two solutions,  $u_i$ ,  $i = 1, 2$  and let  $w = u_1 - u_2$ . Then use the maximum principle.
13. Suppose  $\mathbf{B}$  is a vector field and  $\nabla \times \mathbf{A} = \mathbf{B}$ . Thus  $\mathbf{A}$  is a vector potential for  $\mathbf{B}$ . Show that  $\mathbf{A} + \nabla \phi$  is also a vector potential for  $\mathbf{B}$ . Here  $\phi$  is just a  $C^2$  scalar field. Thus the vector potential is not unique.

### 23.3 The Divergence Theorem

The divergence theorem relates an integral over a set to one on the boundary of the set. It is also called Gauss's theorem.

**Definition 23.3.1** A subset,  $V$  of  $\mathbb{R}^3$  is called cylindrical in the  $x$  direction if it is of the form

$$V = \{(x, y, z) : \phi(y, z) \leq x \leq \psi(y, z) \text{ for } (y, z) \in D\}$$

where  $D$  is a subset of the  $yz$  plane.  $V$  is cylindrical in the  $z$  direction if

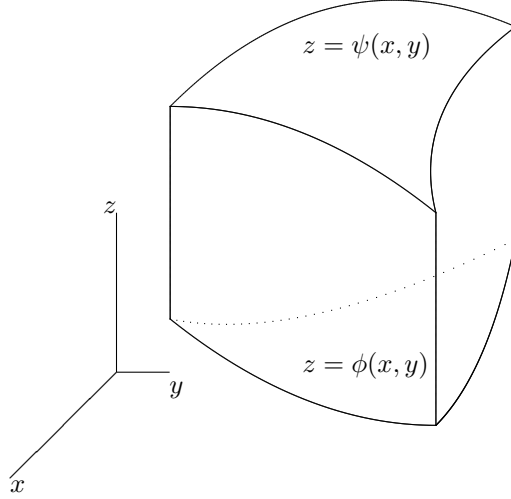
$$V = \{(x, y, z) : \phi(x, y) \leq z \leq \psi(x, y) \text{ for } (x, y) \in D\}$$

where  $D$  is a subset of the  $xy$  plane, and  $V$  is cylindrical in the  $y$  direction if

$$V = \{(x, y, z) : \phi(x, z) \leq y \leq \psi(x, z) \text{ for } (x, z) \in D\}$$

where  $D$  is a subset of the  $xz$  plane. If  $V$  is cylindrical in the  $z$  direction, denote by  $\partial V$  the boundary of  $V$  defined to be the points of the form  $(x, y, \phi(x, y)), (x, y, \psi(x, y))$  for  $(x, y) \in D$ , along with points of the form  $(x, y, z)$  where  $(x, y) \in \partial D$  and  $\phi(x, y) \leq z \leq \psi(x, y)$ . Points on  $\partial D$  are defined to be those for which every open ball contains points which are in  $D$  as well as points which are not in  $D$ . A similar definition holds for  $\partial V$  in the case that  $V$  is cylindrical in one of the other directions.

The following picture illustrates the above definition in the case of  $V$  cylindrical in the  $z$  direction.



Of course, many three dimensional sets are cylindrical in each of the coordinate directions. For example, a ball or a rectangle or a tetrahedron are all cylindrical in each direction. The following lemma allows the exchange of the volume integral of a partial derivative for an area integral in which the derivative is replaced with multiplication by an appropriate component of the unit exterior normal.

**Lemma 23.3.2** Suppose  $V$  is cylindrical in the  $z$  direction and that  $\phi$  and  $\psi$  are the functions in the above definition. Assume  $\phi$  and  $\psi$  are  $C^1$  functions and suppose  $F$  is a  $C^1$  function defined on  $V$ . Also, let  $\mathbf{n} = (n_x, n_y, n_z)$  be the unit exterior normal to  $\partial V$ . Then

$$\iiint_V \frac{\partial F}{\partial z}(x, y, z) dV = \iint_{\partial V} F n_z dA.$$

**Proof:** From the fundamental theorem of calculus,

$$\begin{aligned} \iiint_V \frac{\partial F}{\partial z}(x, y, z) dV &= \iint_D \int_{\phi(x, y)}^{\psi(x, y)} \frac{\partial F}{\partial z}(x, y, z) dz dx dy \\ &= \iint_D [F(x, y, \psi(x, y)) - F(x, y, \phi(x, y))] dx dy \end{aligned} \quad (23.3)$$

Now the unit exterior normal on the top of  $V$ , the surface  $(x, y, \psi(x, y))$  is

$$\frac{1}{\sqrt{\psi_x^2 + \psi_y^2 + 1}} (-\psi_x, -\psi_y, 1).$$

This follows from the observation that the top surface is the level surface,  $z - \psi(x, y) = 0$  and so the gradient of this function of three variables is perpendicular to the level surface. It points in the correct direction because the  $z$  component is positive. Therefore, on the top surface,

$$n_z = \frac{1}{\sqrt{\psi_x^2 + \psi_y^2 + 1}}$$

Similarly, the unit normal to the surface on the bottom is

$$\frac{1}{\sqrt{\phi_x^2 + \phi_y^2 + 1}} (\phi_x, \phi_y, -1)$$

and so on the bottom surface,

$$n_z = \frac{-1}{\sqrt{\phi_x^2 + \phi_y^2 + 1}}$$

Note that here the  $z$  component is negative because since it is the outer normal it must point down. On the lateral surface, the one where  $(x, y) \in \partial D$  and  $z \in [\phi(x, y), \psi(x, y)]$ ,  $n_z = 0$ .

The area element on the top surface is  $dA = \sqrt{\psi_x^2 + \psi_y^2 + 1} dx dy$  while the area element on the bottom surface is  $\sqrt{\phi_x^2 + \phi_y^2 + 1} dx dy$ . Therefore, the last expression in (23.3) is of the form,

$$\begin{aligned} & \int \int_D F(x, y, \psi(x, y)) \overbrace{\frac{1}{\sqrt{\psi_x^2 + \psi_y^2 + 1}}}^{n_z} \overbrace{\sqrt{\psi_x^2 + \psi_y^2 + 1} dx dy}^{dA} + \\ & \int \int_D F(x, y, \phi(x, y)) \left( \overbrace{\frac{-1}{\sqrt{\phi_x^2 + \phi_y^2 + 1}}}^{n_z} \right) \overbrace{\sqrt{\phi_x^2 + \phi_y^2 + 1} dx dy}^{dA} \\ & + \int \int_{\text{Lateral surface}} F n_z dA, \end{aligned}$$

the last term equaling zero because on the lateral surface,  $n_z = 0$ . Therefore, this reduces to  $\int \int_{\partial V} F n_z dA$  as claimed.

The following corollary is entirely similar to the above.

**Corollary 23.3.3** *If  $V$  is cylindrical in the  $y$  direction, then*

$$\int \int \int_V \frac{\partial F}{\partial y} dV = \int \int_{\partial V} F n_y dA$$

*and if  $V$  is cylindrical in the  $x$  direction, then*

$$\int \int \int_V \frac{\partial F}{\partial x} dV = \int \int_{\partial V} F n_x dA$$

With this corollary, here is a proof of the divergence theorem.

**Theorem 23.3.4** *Let  $V$  be cylindrical in each of the coordinate directions and let  $\mathbf{F}$  be a  $C^1$  vector field defined on  $V$ . Then*

$$\int \int \int_V \nabla \cdot \mathbf{F} dV = \int \int_{\partial V} \mathbf{F} \cdot \mathbf{n} dA.$$

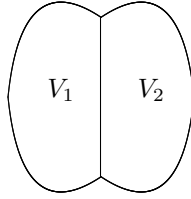


**Proof:** From the above lemma and corollary,

$$\begin{aligned}
 \int \int \int_V \nabla \cdot \mathbf{F} dV &= \int \int \int_V \frac{\partial F_1}{\partial x} + \frac{\partial F_2}{\partial y} + \frac{\partial F_3}{\partial z} dV \\
 &= \int \int_{\partial V} (F_1 n_x + F_2 n_y + F_3 n_z) dA \\
 &= \int \int_{\partial V} \mathbf{F} \cdot \mathbf{n} dA.
 \end{aligned}$$

This proves the theorem.

The divergence theorem holds for much more general regions than this. Suppose for example you have a complicated region which is the union of finitely many disjoint regions of the sort just described which are cylindrical in each of the coordinate directions. Then the volume integral over the union of these would equal the sum of the integrals over the disjoint regions. If the boundaries of two of these regions intersect, then the area integrals will cancel out on the intersection because the unit exterior normals will point in opposite directions. Therefore, the sum of the integrals over the boundaries of these disjoint regions will reduce to an integral over the boundary of the union of these. Hence the divergence theorem will continue to hold. For example, consider the following picture. If the divergence theorem holds for each  $V_i$  in the following picture, then it holds for the union of these two.



General formulations of the divergence theorem involve Hausdorff measures and the Lebesgue integral, a better integral than the old fashioned Riemann integral which has been obsolete now for almost 100 years. When all is said and done, one finds that the conclusion of the divergence theorem is usually true and the theorem can be used with confidence.

This theorem also makes possible a coordinate free definition of the divergence.

**Theorem 23.3.5** Let  $B(\mathbf{x}, \delta)$  be the ball centered at  $\mathbf{x}$  having radius  $\delta$  and let  $\mathbf{F}$  be a  $C^1$  vector field. Then letting  $v(B(\mathbf{x}, \delta))$  denote the volume of  $B(\mathbf{x}, \delta)$  given by

$$\int_{B(\mathbf{x}, \delta)} dV,$$

it follows

$$\operatorname{div} \mathbf{F}(\mathbf{x}) = \lim_{\delta \rightarrow 0^+} \frac{1}{v(B(\mathbf{x}, \delta))} \int_{\partial B(\mathbf{x}, \delta)} \mathbf{F} \cdot \mathbf{n} dA. \quad (23.4)$$

**Proof:** The divergence theorem holds for balls because they are cylindrical in every direction. Therefore,

$$\frac{1}{v(B(\mathbf{x}, \delta))} \int_{\partial B(\mathbf{x}, \delta)} \mathbf{F} \cdot \mathbf{n} dA = \frac{1}{v(B(\mathbf{x}, \delta))} \int_{B(\mathbf{x}, \delta)} \operatorname{div} \mathbf{F}(\mathbf{y}) dV.$$

Therefore, since  $\operatorname{div} \mathbf{F}(\mathbf{x})$  is a constant,

$$\left| \operatorname{div} \mathbf{F}(\mathbf{x}) - \frac{1}{v(B(\mathbf{x}, \delta))} \int_{\partial B(\mathbf{x}, \delta)} \mathbf{F} \cdot \mathbf{n} dA \right|$$

$$\begin{aligned}
&= \left| \operatorname{div} \mathbf{F}(\mathbf{x}) - \frac{1}{v(B(\mathbf{x}, \delta))} \int_{B(\mathbf{x}, \delta)} \operatorname{div} \mathbf{F}(\mathbf{y}) \, dV \right| \\
&= \left| \frac{1}{v(B(\mathbf{x}, \delta))} \int_{B(\mathbf{x}, \delta)} (\operatorname{div} \mathbf{F}(\mathbf{x}) - \operatorname{div} \mathbf{F}(\mathbf{y})) \, dV \right| \\
&\leq \frac{1}{v(B(\mathbf{x}, \delta))} \int_{B(\mathbf{x}, \delta)} |\operatorname{div} \mathbf{F}(\mathbf{x}) - \operatorname{div} \mathbf{F}(\mathbf{y})| \, dV \\
&\leq \frac{1}{v(B(\mathbf{x}, \delta))} \int_{B(\mathbf{x}, \delta)} \frac{\varepsilon}{2} \, dV < \varepsilon
\end{aligned}$$

whenever  $\varepsilon$  is small enough due to the continuity of  $\operatorname{div} \mathbf{F}$ . Since  $\varepsilon$  is arbitrary, this shows (23.4).

How is this definition independent of coordinates? It only involves geometrical notions of volume and dot product. This is why. Imagine rotating the coordinate axes, keeping all distances the same and expressing everything in terms of the new coordinates. The divergence would still have the same value because of this theorem.

## 23.4 Exercises

1. To prove the divergence theorem, it was shown first that the spacial partial derivative in the volume integral could be exchanged for multiplication by an appropriate component of the exterior normal. This problem starts with the divergence theorem and goes the other direction. Assuming the divergence theorem, holds for a region,  $V$ , show that  $\int_{\partial V} \mathbf{n} u \, dA = \int_V \nabla u \, dV$ . Note this implies  $\int_V \frac{\partial u}{\partial x} \, dV = \int_{\partial V} n_1 u \, dA$ .
2. Let  $V$  be such that the divergence theorem holds. Show that  $\int_V \nabla \cdot (u \nabla v) \, dV = \int_{\partial V} u \frac{\partial v}{\partial \mathbf{n}} \, dA$  where  $\mathbf{n}$  is the exterior normal and  $\frac{\partial v}{\partial \mathbf{n}}$  denotes the directional derivative of  $v$  in the direction  $\mathbf{n}$ .
3. Let  $V$  be such that the divergence theorem holds. Show that  $\int \int_V (v \nabla^2 u - u \nabla^2 v) \, dV = \int \int_{\partial V} (v \frac{\partial u}{\partial \mathbf{n}} - u \frac{\partial v}{\partial \mathbf{n}}) \, dA$  where  $\mathbf{n}$  is the exterior normal and  $\frac{\partial u}{\partial \mathbf{n}}$  is defined in Problem 2.
4. Let  $V$  be a ball and suppose  $\nabla^2 u = f$  in  $V$  while  $u = g$  on  $\partial V$ . Show there is at most one solution to this boundary value problem which is  $C^2$  in  $V$  and continuous on  $V$  with its boundary. **Hint:** You might consider  $w = u - v$  where  $u$  and  $v$  are solutions to the problem. Then use the result of Problem 2 and the identity

$$w \nabla^2 w = \nabla \cdot (w \nabla w) - \nabla w \cdot \nabla w$$

to conclude  $\nabla w = 0$ . Then show this implies  $w$  must be a constant by considering  $h(t) = w(t\mathbf{x} + (1-t)\mathbf{y})$  and showing  $h$  is a constant. Alternatively, you might consider the maximum principle.

5. Show that  $\int_{\partial V} \nabla \times \mathbf{v} \cdot \mathbf{n} \, dA = 0$  where  $V$  is a region for which the divergence theorem holds and  $\mathbf{v}$  is a  $C^2$  vector field.
6. Let  $\mathbf{F}(x, y, z) = (x, y, z)$  be a vector field in  $\mathbb{R}^3$  and let  $V$  be a three dimensional shape and let  $\mathbf{n} = (n_1, n_2, n_3)$ . Show  $\int_{\partial V} (xn_1 + yn_2 + zn_3) \, dA = 3 \times \text{volume of } V$ .
7. Does the divergence theorem hold for higher dimensions? If so, explain why it does. How about two dimensions?
8. Let  $\mathbf{F} = x\mathbf{i} + y\mathbf{j} + z\mathbf{k}$  and let  $V$  denote the tetrahedron formed by the planes,  $x = 0, y = 0, z = 0$ , and  $\frac{1}{3}x + \frac{1}{3}y + \frac{1}{3}z = 1$ . Verify the divergence theorem for this example.

## 23.5 Some Applications Of The Divergence Theorem

### 23.5.1 Hydrostatic Pressure

Imagine a fluid which does not move which is acted on by an acceleration,  $\mathbf{g}$ . Of course the acceleration is usually the acceleration of gravity. Also let the density of the fluid be  $\rho$ , a function of position. What can be said about the pressure,  $p$ , in the fluid? Let  $B(\mathbf{x}, \varepsilon)$  be a small ball centered at the point,  $\mathbf{x}$ . Then the force the fluid exerts on this ball would equal

$$-\int_{\partial B(\mathbf{x}, \varepsilon)} p \mathbf{n} dA.$$

Here  $\mathbf{n}$  is the unit exterior normal at a small piece of  $\partial B(\mathbf{x}, \varepsilon)$  having area  $dA$ . By the divergence theorem, (see Problem 1 on Page 602) this integral equals

$$-\int_{B(\mathbf{x}, \varepsilon)} \nabla p dV.$$

Also the force acting on this small ball of fluid is

$$\int_{B(\mathbf{x}, \varepsilon)} \rho \mathbf{g} dV.$$

Since it is given that the fluid does not move, the sum of these forces must equal zero. Thus

$$\int_{B(\mathbf{x}, \varepsilon)} \rho \mathbf{g} dV = \int_{B(\mathbf{x}, \varepsilon)} \nabla p dV.$$

Since this must hold for any ball in the fluid of any radius, it must be that

$$\nabla p = \rho \mathbf{g}. \quad (23.5)$$

Earlier the topic of force on a dam was discussed where it was asserted that the pressure at depth  $z$  was equal to  $62.5z$ . This is easy to see from (23.5). In this case,  $\mathbf{g} = g\mathbf{k}$  where  $g = 32$  feet/sec<sup>2</sup>. The weight of a cubic foot of water is 62.5 pounds. Therefore, the mass in slugs of this water is  $62.5/32$ . Since it is a cubic foot, this is also the density of the water in slugs per cubic foot. Also, it is normally assumed that water is incompressible<sup>1</sup>. Therefore, this is the mass of water at any depth. Therefore,

$$\frac{\partial p}{\partial x} \mathbf{i} + \frac{\partial p}{\partial y} \mathbf{j} + \frac{\partial p}{\partial z} \mathbf{k} = \frac{62.5}{32} \times 32 \mathbf{k}.$$

and so  $p$  does not depend on  $x$  and  $y$  and is only a function of  $z$ . It follows  $p(0) = 0$ , and  $p'(z) = 62.5$ . Therefore,  $p(x, y, z) = 62.5z$ . This establishes the earlier claim. This is interesting but (23.5) is more interesting because it does not require  $\rho$  to be constant.

### 23.5.2 Archimedes Law Of Buoyancy

Archimedes principle states that when a solid body is immersed in a fluid the net force acting on the body by the fluid is directly up and equals the total weight of the fluid displaced.

Denote the set of points in three dimensions occupied by the body as  $V$ . Then for  $dA$  an increment of area on the surface of this body, the force acting on this increment of area

---

<sup>1</sup>There is no such thing as an incompressible fluid but this doesn't stop people from making this assumption.

would equal  $-p d\mathbf{A}\mathbf{n}$  where  $\mathbf{n}$  is the exterior unit normal. Therefore, since the fluid does not move,

$$\int_{\partial V} -p\mathbf{n} dA = \int_V -\nabla p dV = \int_V \rho g dV \mathbf{k}$$

Which equals the total weight of the displaced fluid and you note the force is directed upward as claimed. Here  $\rho$  is the density and (23.5) is being used. There is an interesting point in the above explanation. Why does the second equation hold? Imagine that  $V$  were filled with fluid. Then the equation follows from (23.5) because in this equation  $\mathbf{g} = -g\mathbf{k}$ .

### 23.5.3 Equations Of Heat And Diffusion

Let  $\mathbf{x}$  be a point in three dimensional space and let  $(x_1, x_2, x_3)$  be Cartesian coordinates of this point. Let there be a three dimensional body having density,  $\rho = \rho(\mathbf{x}, t)$ .

The heat flux,  $\mathbf{J}$ , in the body is defined as a vector which has the following property.

$$\text{Rate at which heat crosses } S = \int_S \mathbf{J} \cdot \mathbf{n} dA$$

where  $\mathbf{n}$  is the unit normal in the desired direction. Thus if  $V$  is a three dimensional body,

$$\text{Rate at which heat leaves } V = \int_{\partial V} \mathbf{J} \cdot \mathbf{n} dA$$

where  $\mathbf{n}$  is the unit exterior normal.

Fourier's law of heat conduction states that the heat flux,  $\mathbf{J}$  satisfies  $\mathbf{J} = -k\nabla(u)$  where  $u$  is the temperature and  $k = k(u, \mathbf{x}, t)$  is called the coefficient of thermal conductivity. This changes depending on the material. It also can be shown by experiment to change with temperature. This equation for the heat flux states that the heat flows from hot places toward colder places in the direction of greatest rate of decrease in temperature. Let  $c(\mathbf{x}, t)$  denote the specific heat of the material in the body. This means the amount of heat within  $V$  is given by the formula  $\int \int_V \rho(\mathbf{x}, t) c(\mathbf{x}, t) u(\mathbf{x}, t) dV$ . Suppose also there are sources for the heat within the material given by  $f(\mathbf{x}, u, t)$ . If  $f$  is positive, the heat is increasing while if  $f$  is negative the heat is decreasing. For example such sources could result from a chemical reaction taking place. Then the divergence theorem can be used to verify the following equation for  $u$ . Such an equation is called a reaction diffusion equation.

$$\frac{\partial}{\partial t} (\rho(\mathbf{x}, t) c(\mathbf{x}, t) u(\mathbf{x}, t)) = \nabla \cdot (k(u, \mathbf{x}, t) \nabla u(\mathbf{x}, t)) + f(\mathbf{x}, u, t). \quad (23.6)$$

Take an arbitrary  $V$  for which the divergence theorem holds. Then the time rate of change of the heat in  $V$  is

$$\frac{d}{dt} \int_V \rho(\mathbf{x}, t) c(\mathbf{x}, t) u(\mathbf{x}, t) dV = \int_V \frac{\partial (\rho(\mathbf{x}, t) c(\mathbf{x}, t) u(\mathbf{x}, t))}{\partial t} dV$$

where, as in the preceding example, this is a physical derivation so the consideration of hard mathematics is not necessary. Therefore, from the Fourier law of heat conduction,  $\frac{d}{dt} \int_V \rho(\mathbf{x}, t) c(\mathbf{x}, t) u(\mathbf{x}, t) dV =$

$$\begin{aligned} \int_V \frac{\partial (\rho(\mathbf{x}, t) c(\mathbf{x}, t) u(\mathbf{x}, t))}{\partial t} dV &= \overbrace{\int_{\partial V} -\mathbf{J} \cdot \mathbf{n} dA}^{\text{rate at which heat enters}} + \int_V f(\mathbf{x}, u, t) dV \\ &= \int_{\partial V} k \nabla(u) \cdot \mathbf{n} dA + \int_V f(\mathbf{x}, u, t) dV = \int \int \int_V (\nabla \cdot (k \nabla(u)) + f) dV. \end{aligned}$$

Since this holds for every sample volume,  $V$  it must be the case that the above reaction diffusion equation, (23.6) holds. Note that more interesting equations can be obtained by letting more of the quantities in the equation depend on temperature. However, the above is a fairly hard equation and people usually assume the coefficient of thermal conductivity depends only on  $\mathbf{x}$  and that the reaction term,  $f$  depends only on  $\mathbf{x}$  and  $t$  and that  $\rho$  and  $c$  are constant. Then it reduces to the much easier equation,

$$\frac{\partial}{\partial t} u(\mathbf{x}, t) = \frac{1}{\rho c} \nabla \cdot (k(\mathbf{x}) \nabla u(\mathbf{x}, t)) + f(\mathbf{x}, t). \quad (23.7)$$

This is often referred to as the heat equation. Sometimes there are modifications of this in which  $k$  is not just a scalar but a matrix to account for different heat flow properties in different directions. However, they are not much harder than the above. The major mathematical difficulties result from allowing  $k$  to depend on temperature.

It is known that the heat equation is not correct even if the thermal conductivity did not depend on  $u$  because it implies infinite speed of propagation of heat. However, this does not prevent people from using it.

### 23.5.4 Balance Of Mass

Let  $\mathbf{y}$  be a point in three dimensional space and let  $(y_1, y_2, y_3)$  be Cartesian coordinates of this point. Let  $V$  be a region in three dimensional space and suppose a fluid having density,  $\rho(\mathbf{y}, t)$  and velocity,  $\mathbf{v}(\mathbf{y}, t)$  is flowing through this region. Then the mass of fluid leaving  $V$  per unit time is given by the area integral,  $\int_{\partial V} \rho(\mathbf{y}, t) \mathbf{v}(\mathbf{y}, t) \cdot \mathbf{n} dA$  while the total mass of the fluid enclosed in  $V$  at a given time is  $\int_V \rho(\mathbf{y}, t) dV$ . Also suppose mass originates at the rate  $f(\mathbf{y}, t)$  per cubic unit per unit time within this fluid. Then the conclusion which can be drawn through the use of the divergence theorem is the following fundamental equation known as the mass balance equation.

$$\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho \mathbf{v}) = f(\mathbf{y}, t) \quad (23.8)$$

To see this is so, take an arbitrary  $V$  for which the divergence theorem holds. Then the time rate of change of the mass in  $V$  is

$$\frac{\partial}{\partial t} \int_V \rho(\mathbf{y}, t) dV = \int_V \frac{\partial \rho(\mathbf{y}, t)}{\partial t} dV$$

where the derivative was taken under the integral sign with respect to  $t$ . (This is a physical derivation and therefore, it is not necessary to fuss with the hard mathematics related to the change of limit operations. You should expect this to be true under fairly general conditions because the integral is a sort of sum and the derivative of a sum is the sum of the derivatives.) Therefore, the rate of change of mass,  $\frac{\partial}{\partial t} \int_V \rho(\mathbf{y}, t) dV$ , equals

$$\begin{aligned} \int_V \frac{\partial \rho(\mathbf{y}, t)}{\partial t} dV &= \overbrace{\int_{\partial V} \rho(\mathbf{y}, t) \mathbf{v}(\mathbf{y}, t) \cdot \mathbf{n} dA}^{\text{rate at which mass leaves}} + \int_V f(\mathbf{y}, t) dV \\ &= \int_V (\nabla \cdot (\rho(\mathbf{y}, t) \mathbf{v}(\mathbf{y}, t)) + f(\mathbf{y}, t)) dV. \end{aligned}$$

Since this holds for every sample volume,  $V$  it must be the case that the equation of continuity holds. Again, there are interesting mathematical questions here which can be explored but since it is a physical derivation, it is not necessary to dwell too much on them.

If all the functions involved are continuous, it is certainly true but it is true under far more general conditions than that.

Also note this equation applies to many situations and  $f$  might depend on more than just  $\mathbf{y}$  and  $t$ . In particular,  $f$  might depend also on temperature and the density,  $\rho$ . This would be the case for example if you were considering the mass of some chemical and  $f$  represented a chemical reaction. Mass balance is a general sort of equation valid in many contexts.

### 23.5.5 Balance Of Momentum

This example is a little more substantial than the above. It concerns the balance of momentum for a continuum. To see a full description of all the physics involved, you should consult a book on continuum mechanics. The situation is of a material in three dimensions and it deforms and moves about in three dimensions. This means this material is not a rigid body. Let  $B_0$  denote an open set identifying a chunk of this material at time  $t = 0$  and let  $B_t$  be an open set which identifies the same chunk of material at time  $t > 0$ .

Let  $\mathbf{y}(t, \mathbf{x}) = (y_1(t, \mathbf{x}), y_2(t, \mathbf{x}), y_3(t, \mathbf{x}))$  denote the position with respect to Cartesian coordinates at time  $t$  of the point whose position at time  $t = 0$  is  $\mathbf{x} = (x_1, x_2, x_3)$ . The coordinates,  $\mathbf{x}$  are sometimes called the reference coordinates and sometimes the material coordinates and sometimes the Lagrangian coordinates. The coordinates,  $\mathbf{y}$  are called the Eulerian coordinates or sometimes the spacial coordinates and the function,  $(t, \mathbf{x}) \rightarrow \mathbf{y}(t, \mathbf{x})$  is called the motion. Thus

$$\mathbf{y}(0, \mathbf{x}) = \mathbf{x}. \quad (23.9)$$

The derivative,

$$D_2 \mathbf{y}(t, \mathbf{x})$$

is called the deformation gradient. Recall the notation means you fix  $t$  and consider the function,  $\mathbf{x} \rightarrow \mathbf{y}(t, \mathbf{x})$ , taking its derivative. Since it is a linear transformation, it is represented by the usual matrix, whose  $ij^{th}$  entry is given by

$$F_{ij}(\mathbf{x}) = \frac{\partial y_i(t, \mathbf{x})}{\partial x_j}.$$

Let  $\rho(t, \mathbf{y})$  denote the density of the material at time  $t$  at the point,  $\mathbf{y}$  and let  $\rho_0(\mathbf{x})$  denote the density of the material at the point,  $\mathbf{x}$ . Thus  $\rho_0(\mathbf{x}) = \rho(0, \mathbf{x}) = \rho(0, \mathbf{y}(0, \mathbf{x}))$ . The first task is to consider the relationship between  $\rho(t, \mathbf{y})$  and  $\rho_0(\mathbf{x})$ .

**Lemma 23.5.1**  $\rho_0(\mathbf{x}) = \rho(t, \mathbf{y}(t, \mathbf{x})) \det(F)$  and in any reasonable physical motion,  $\det(F) > 0$ .

**Proof:** Let  $V_0$  represent a small chunk of material at  $t = 0$  and let  $V_t$  represent the same chunk of material at time  $t$ . I will be a little sloppy and refer to  $V_0$  as the small chunk of material at time  $t = 0$  and  $V_t$  as the chunk of material at time  $t$  rather than an open set representing the chunk of material. Then by the change of variables formula for multiple integrals,

$$\int_{V_t} dV = \int_{V_0} |\det(F)| dV.$$

If  $\det(F) = 0$  for some  $t$  the above formula shows that the chunk of material went from positive volume to zero volume and this is not physically possible. Therefore, it is impossible that  $\det(F)$  can equal zero. However, at  $t = 0$ ,  $F = I$ , the identity because of (23.9). Therefore,  $\det(F) = 1$  at  $t = 0$  and if it is assumed  $t \rightarrow \det(F)$  is continuous it follows by the intermediate value theorem that  $\det(F) > 0$  for all  $t$ . Of course it is not known for

sure this function is continuous but the above shows why it is at least reasonable to expect  $\det(F) > 0$ .

Now using the change of variables formula,

$$\begin{aligned} \text{mass of } V_t &= \int_{V_t} \rho(t, \mathbf{y}) \, dV = \int_{V_0} \rho(t, \mathbf{y}(t, \mathbf{x})) \det(F) \, dV \\ &= \text{mass of } V_0 = \int_{V_0} \rho_0(\mathbf{x}) \, dV. \end{aligned}$$

Since  $V_0$  is arbitrary, it follows  $\rho_0(\mathbf{x}) = \rho(t, \mathbf{y}(t, \mathbf{x})) \det(F)$  as claimed. Note this shows that  $\det(F)$  is a magnification factor for the density.

Now consider a small chunk of material,  $B_t$  at time  $t$  which corresponds to  $B_0$  at time  $t = 0$ . The total linear momentum of this material at time  $t$  is

$$\int_{B_t} \rho(t, \mathbf{y}) \mathbf{v}(t, \mathbf{y}) \, dV$$

where  $\mathbf{v}$  is the velocity. By Newton's second law, the time rate of change of this linear momentum should equal the total force acting on the chunk of material. In the following derivation,  $dV(\mathbf{y})$  will indicate the integration is taking place with respect to the variable,  $\mathbf{y}$ . By Lemma 23.5.1 and the change of variables formula for multiple integrals

$$\begin{aligned} \frac{d}{dt} \left( \int_{B_t} \rho(t, \mathbf{y}) \mathbf{v}(t, \mathbf{y}) \, dV(\mathbf{y}) \right) &= \frac{d}{dt} \left( \int_{B_0} \rho(t, \mathbf{y}(t, \mathbf{x})) \mathbf{v}(t, \mathbf{y}(t, \mathbf{x})) \det(F) \, dV(\mathbf{x}) \right) \\ &= \frac{d}{dt} \left( \int_{B_0} \rho_0(\mathbf{x}) \mathbf{v}(t, \mathbf{y}(t, \mathbf{x})) \, dV(\mathbf{x}) \right) \\ &= \int_{B_0} \rho_0(\mathbf{x}) \left[ \frac{\partial \mathbf{v}}{\partial t} + \frac{\partial \mathbf{v}}{\partial y_i} \frac{\partial y_i}{\partial t} \right] dV(\mathbf{x}) \\ &= \int_{B_t} \rho(t, \mathbf{y}) \det(F) \left[ \frac{\partial \mathbf{v}}{\partial t} + \frac{\partial \mathbf{v}}{\partial y_i} \frac{\partial y_i}{\partial t} \right] \frac{1}{\det(F)} dV(\mathbf{y}) \\ &= \int_{B_t} \rho(t, \mathbf{y}) \left[ \frac{\partial \mathbf{v}}{\partial t} + \frac{\partial \mathbf{v}}{\partial y_i} \frac{\partial y_i}{\partial t} \right] dV(\mathbf{y}). \end{aligned}$$

Having taken the derivative of the total momentum, it is time to consider the total force acting on the chunk of material.

The force comes from two sources, a body force,  $\mathbf{b}$  and a force which act on the boundary of the chunk of material called a traction force. Typically, the body force is something like gravity in which case,  $\mathbf{b} = -g\rho\mathbf{k}$ , assuming the Cartesian coordinate system has been chosen in the usual manner. The traction force is of the form

$$\int_{\partial B_t} \mathbf{s}(t, \mathbf{y}, \mathbf{n}) \, dA$$

where  $\mathbf{n}$  is the unit exterior normal. Thus the traction force depends on position, time, and the orientation of the boundary of  $B_t$ . Cauchy showed the existence of a linear transformation,  $T(t, \mathbf{y})$  such that  $T(t, \mathbf{y}) \mathbf{n} = \mathbf{s}(t, \mathbf{y}, \mathbf{n})$ . It follows there is a matrix,  $T_{ij}(t, \mathbf{y})$  such that the  $i^{th}$  component of  $\mathbf{s}$  is given by  $\mathbf{s}_i(t, \mathbf{y}, \mathbf{n}) = T_{ij}(t, \mathbf{y}) n_j$ . Cauchy also showed this matrix is symmetric,  $T_{ij} = T_{ji}$ . It is called the Cauchy stress. Using Newton's second law to equate the time derivative of the total linear momentum with the applied forces and using the usual repeated index summation convention,

$$\int_{B_t} \rho(t, \mathbf{y}) \left[ \frac{\partial \mathbf{v}}{\partial t} + \frac{\partial \mathbf{v}}{\partial y_i} \frac{\partial y_i}{\partial t} \right] dV(\mathbf{y}) = \int_{B_t} \mathbf{b}(t, \mathbf{y}) \, dV(\mathbf{y}) + \int_{\partial B_t} T_{ij}(t, \mathbf{y}) n_j \, dA.$$

Here is where the divergence theorem is used. In the last integral, the multiplication by  $n_j$  is exchanged for the  $j^{th}$  partial derivative and an integral over  $B_t$ . Thus

$$\int_{B_t} \rho(t, \mathbf{y}) \left[ \frac{\partial \mathbf{v}}{\partial t} + \frac{\partial \mathbf{v}}{\partial y_i} \frac{\partial y_i}{\partial t} \right] dV(\mathbf{y}) = \int_{B_t} \mathbf{b}(t, \mathbf{y}) dV(\mathbf{y}) + \int_{B_t} \frac{\partial (T_{ij}(t, \mathbf{y}))}{\partial y_j} dV(\mathbf{y}).$$

Since  $B_t$  was arbitrary, it follows

$$\begin{aligned} \rho(t, \mathbf{y}) \left[ \frac{\partial \mathbf{v}}{\partial t} + \frac{\partial \mathbf{v}}{\partial y_i} \frac{\partial y_i}{\partial t} \right] &= \mathbf{b}(t, \mathbf{y}) + \frac{\partial (T_{ij}(t, \mathbf{y}))}{\partial y_j} \\ &\equiv \mathbf{b}(t, \mathbf{y}) + \operatorname{div}(T) \end{aligned}$$

where here  $\operatorname{div} T$  is a vector whose  $i^{th}$  component is given by

$$(\operatorname{div} T)_i = \frac{\partial T_{ij}}{\partial y_j}.$$

The term,  $\frac{\partial \mathbf{v}}{\partial t} + \frac{\partial \mathbf{v}}{\partial y_i} \frac{\partial y_i}{\partial t}$ , is the total derivative with respect to  $t$  of the velocity  $\mathbf{v}$ . Thus you might see this written as

$$\rho \dot{\mathbf{v}} = \mathbf{b} + \operatorname{div}(T).$$

The above formulation of the balance of momentum involves the spatial coordinates,  $\mathbf{y}$  but people also like to formulate momentum balance in terms of the material coordinates,  $\mathbf{x}$ . Of course this changes everything.

The momentum in terms of the material coordinates is

$$\int_{B_0} \rho_0(\mathbf{x}) \mathbf{v}(t, \mathbf{x}) dV$$

and so, since  $\mathbf{x}$  does not depend on  $t$ ,

$$\frac{d}{dt} \left( \int_{B_0} \rho_0(\mathbf{x}) \mathbf{v}(t, \mathbf{x}) dV \right) = \int_{B_0} \rho_0(\mathbf{x}) \mathbf{v}_t(t, \mathbf{x}) dV.$$

As indicated earlier, this is a physical derivation and so the mathematical questions related to interchange of limit operations are ignored. This must equal the total applied force. Thus

$$\int_{B_0} \rho_0(\mathbf{x}) \mathbf{v}_t(t, \mathbf{x}) dV = \int_{B_0} \mathbf{b}_0(t, \mathbf{x}) dV + \int_{\partial B_t} T_{ij} n_j dA, \quad (23.10)$$

the first term on the right being the contribution of the body force given per unit volume in the material coordinates and the last term being the traction force discussed earlier. The task is to write this last integral as one over  $\partial B_0$ . For  $\mathbf{y} \in \partial B_t$  there is a unit outer normal,  $\mathbf{n}$ . Here  $\mathbf{y} = \mathbf{y}(t, \mathbf{x})$  for  $\mathbf{x} \in \partial B_0$ . Then define  $\mathbf{N}$  to be the unit outer normal to  $B_0$  at the point,  $\mathbf{x}$ . Near the point  $\mathbf{y} \in \partial B_t$  the surface,  $\partial B_t$  is given parametrically in the form  $\mathbf{y} = \mathbf{y}(s, t)$  for  $(s, t) \in D \subseteq \mathbb{R}^2$  and it can be assumed the unit normal to  $\partial B_t$  near this point is

$$\mathbf{n} = \frac{\mathbf{y}_s(s, t) \times \mathbf{y}_t(s, t)}{|\mathbf{y}_s(s, t) \times \mathbf{y}_t(s, t)|}$$

with the area element given by  $|\mathbf{y}_s(s, t) \times \mathbf{y}_t(s, t)| ds dt$ . This is true for  $\mathbf{y} \in P_t \subseteq \partial B_t$ , a small piece of  $\partial B_t$ . Therefore, the last integral in (23.10) is the sum of integrals over small pieces of the form

$$\int_{P_t} T_{ij} n_j dA \quad (23.11)$$



where  $P_t$  is parametrized by  $\mathbf{y}(s, t)$ ,  $(s, t) \in D$ . Thus the integral in (23.11) is of the form

$$\int_D T_{ij}(\mathbf{y}(s, t)) (\mathbf{y}_s(s, t) \times \mathbf{y}_t(s, t))_j ds dt.$$

By the chain rule this equals

$$\int_D T_{ij}(\mathbf{y}(s, t)) \left( \frac{\partial \mathbf{y}}{\partial x_\alpha} \frac{\partial x_\alpha}{\partial s} \times \frac{\partial \mathbf{y}}{\partial x_\beta} \frac{\partial x_\beta}{\partial t} \right)_j ds dt.$$

Remember  $\mathbf{y} = \mathbf{y}(t, \mathbf{x})$  and it is always assumed the mapping  $\mathbf{x} \rightarrow \mathbf{y}(t, \mathbf{x})$  is one to one and so, since on the surface  $\partial B_t$  near  $\mathbf{y}$ , the points are functions of  $(s, t)$ , it follows  $\mathbf{x}$  is also a function of  $(s, t)$ . Now by the properties of the cross product, this last integral equals

$$\int_D T_{ij}(\mathbf{x}(s, t)) \frac{\partial x_\alpha}{\partial s} \frac{\partial x_\beta}{\partial t} \left( \frac{\partial \mathbf{y}}{\partial x_\alpha} \times \frac{\partial \mathbf{y}}{\partial x_\beta} \right)_j ds dt \quad (23.12)$$

where here  $\mathbf{x}(s, t)$  is the point of  $\partial B_0$  which corresponds with  $\mathbf{y}(s, t) \in \partial B_t$ . Thus  $T_{ij}(\mathbf{x}(s, t)) = T_{ij}(\mathbf{y}(s, t))$ . (Perhaps this is a slight abuse of notation because  $T_{ij}$  is defined on  $\partial B_t$ , not on  $\partial B_0$ , but it avoids introducing extra symbols.) Next (23.12) equals

$$\begin{aligned} & \int_D T_{ij}(\mathbf{x}(s, t)) \frac{\partial x_\alpha}{\partial s} \frac{\partial x_\beta}{\partial t} \varepsilon_{jab} \frac{\partial y_a}{\partial x_\alpha} \frac{\partial y_b}{\partial x_\beta} ds dt \\ &= \int_D T_{ij}(\mathbf{x}(s, t)) \frac{\partial x_\alpha}{\partial s} \frac{\partial x_\beta}{\partial t} \varepsilon_{cab} \delta_{jc} \frac{\partial y_a}{\partial x_\alpha} \frac{\partial y_b}{\partial x_\beta} ds dt \\ &= \int_D T_{ij}(\mathbf{x}(s, t)) \frac{\partial x_\alpha}{\partial s} \frac{\partial x_\beta}{\partial t} \varepsilon_{cab} \overbrace{\frac{\partial y_c}{\partial x_p} \frac{\partial x_p}{\partial y_j}}^{=\delta_{jc}} \frac{\partial y_a}{\partial x_\alpha} \frac{\partial y_b}{\partial x_\beta} ds dt \\ &= \int_D T_{ij}(\mathbf{x}(s, t)) \frac{\partial x_\alpha}{\partial s} \frac{\partial x_\beta}{\partial t} \frac{\partial x_p}{\partial y_j} \overbrace{\varepsilon_{cab} \frac{\partial y_c}{\partial x_p} \frac{\partial y_a}{\partial x_\alpha} \frac{\partial y_b}{\partial x_\beta}}^{=\varepsilon_{p\alpha\beta} \det(F)} ds dt \\ &= \int_D (\det F) T_{ij}(\mathbf{x}(s, t)) \varepsilon_{p\alpha\beta} \frac{\partial x_\alpha}{\partial s} \frac{\partial x_\beta}{\partial t} \frac{\partial x_p}{\partial y_j} ds dt. \end{aligned}$$

Now  $\frac{\partial x_p}{\partial y_j} = F_{pj}^{-1}$  and also

$$\varepsilon_{p\alpha\beta} \frac{\partial x_\alpha}{\partial s} \frac{\partial x_\beta}{\partial t} = (\mathbf{x}_s \times \mathbf{x}_t)_p$$

so the result just obtained is of the form

$$\begin{aligned} & \int_D (\det F) F_{pj}^{-1} T_{ij}(\mathbf{x}(s, t)) (\mathbf{x}_s \times \mathbf{x}_t)_p ds dt = \\ & \int_D (\det F) T_{ij}(\mathbf{x}(s, t)) (F^{-T})_{jp} (\mathbf{x}_s \times \mathbf{x}_t)_p ds dt. \end{aligned}$$

This has transformed the integral over  $P_t$  to one over  $P_0$ , the part of  $\partial B_0$  which corresponds with  $P_t$ . Thus the last integral is of the form

$$\int_{P_0} \det(F) (F^{-T} T)_{ip} N_p dA$$

Summing these up over the pieces of  $\partial B_t$  and  $\partial B_0$  yields the last integral in (23.10) equals

$$\int_{\partial B_0} \det(F) (F^{-T}T)_{ip} N_p dA$$

and so the balance of momentum in terms of the material coordinates becomes

$$\int_{B_0} \rho_0(\mathbf{x}) \mathbf{v}_t(t, \mathbf{x}) dV = \int_{B_0} \mathbf{b}_0(t, \mathbf{x}) dV + \int_{\partial B_0} \det(F) (F^{-T}T)_{ip} N_p dA$$

The matrix,  $\det(F) (F^{-T}T)_{ip}$  is called the Piola Kirchhoff stress,  $S$ . An application of the divergence theorem yields

$$\int_{B_0} \rho_0(\mathbf{x}) \mathbf{v}_t(t, \mathbf{x}) dV = \int_{B_0} \mathbf{b}_0(t, \mathbf{x}) dV + \int_{B_0} \frac{\partial (\det(F) (F^{-T}T)_{ip})}{\partial x_p} dV.$$

Since  $B_0$  is arbitrary, a balance law for momentum in terms of the material coordinates is obtained

$$\begin{aligned} \rho_0(\mathbf{x}) \mathbf{v}_t(t, \mathbf{x}) &= \mathbf{b}_0(t, \mathbf{x}) + \frac{\partial (\det(F) (F^{-T}T)_{ip})}{\partial x_p} \\ &= \mathbf{b}_0(t, \mathbf{x}) + \operatorname{div} (\det(F) (F^{-T}T)) \\ &= \mathbf{b}_0(t, \mathbf{x}) + \operatorname{div} S. \end{aligned} \quad (23.13)$$

The main purpose of this presentation is to show how the divergence theorem is used in a significant way to obtain balance laws and to indicate a very interesting direction for further study. To continue, one needs to specify  $T$  or  $S$  as an appropriate function of things related to the motion,  $\mathbf{y}$ . Often the thing related to the motion is something called the strain and such relationships between the stress and the strain are known as constitutive laws. The proper formulation of constitutive laws involves more physical considerations such as frame indifference in which it is required the response of the system cannot depend on the manner in which the Cartesian coordinate system was chosen. There are also many other physical properties which can be included and which require a certain form for the constitutive equations. These considerations are outside the scope of this book and require a considerable amount of linear algebra.

There are also balance laws for energy which you may study later but these are more problematic than the balance laws for mass and momentum. However, the divergence theorem is used in these also.

### 23.5.6 The Wave Equation

As an example of how the balance law of momentum is used to obtain an important equation of mathematical physics, suppose  $S = kF$  where  $k$  is a constant and  $F$  is the deformation gradient and let  $\mathbf{u} \equiv \mathbf{y} - \mathbf{x}$ . Thus  $\mathbf{u}$  is the displacement. Then from (23.13) you can verify the following holds.

$$\rho_0(\mathbf{x}) \mathbf{u}_{tt}(t, \mathbf{x}) = \mathbf{b}_0(t, \mathbf{x}) + k\Delta \mathbf{u}(t, \mathbf{x}) \quad (23.14)$$

In the case where  $\rho_0$  is a constant and  $\mathbf{b}_0 = 0$ , this yields

$$\mathbf{u}_{tt} - c\Delta \mathbf{u} = \mathbf{0}.$$

The wave equation is  $u_{tt} - c\Delta u = 0$  and so the above gives three wave equations, one for each component.

### 23.5.7 A Negative Observation

Many of the above applications of the divergence theorem are based on the assumption that matter is continuously distributed in a way that the above arguments are correct. In other words, a continuum. However, there is no such thing as a continuum. It has been known for some time now that matter is composed of atoms. It is not continuously distributed through some region of space as it is in the above. Apologists for this contradiction with reality sometimes say to consider enough of the material in question that it is reasonable to think of it as a continuum. This mystical reasoning is then violated as soon as they go from the integral form of the balance laws to the differential equations expressing the traditional formulation of these laws. See Problem 1 below, for example. However, these laws continue to be used and seem to lead to useful physical models which have value in predicting the behavior of physical systems. This is what justifies their use, not any fundamental truth.

## 23.6 Exercises

1. Suppose  $f : U \rightarrow \mathbb{R}$  is continuous where  $U$  is some open set and for all  $B \subseteq U$  where  $B$  is a ball,  $\int_B f(\mathbf{x}) dV = 0$ . Show this implies  $f(\mathbf{x}) = 0$  for all  $\mathbf{x} \in U$ .
2. Let  $U$  denote the box centered at  $(0, 0, 0)$  with sides parallel to the coordinate planes which has width 4, length 2 and height 3. Find the flux integral  $\int_{\partial U} \mathbf{F} \cdot \mathbf{n} dS$  where  $\mathbf{F} = \langle x + 3, 2y, 3z \rangle$ . **Hint:** If you like, you might want to use the divergence theorem.
3. Verify (23.14) from (23.13) and the assumption that  $S = kF$ .
4. Fick's law for diffusion states the flux of a diffusing species,  $\mathbf{J}$  is proportional to the gradient of the concentration,  $c$ . Write this law getting the sign right for the constant of proportionality and derive an equation similar to the heat equation for the concentration,  $c$ . Typically,  $c$  is the concentration of some sort of pollutant or a chemical.
5. Show that if  $u_k, k = 1, 2, \dots, n$  each satisfies (23.7) then for any choice of constants,  $c_1, \dots, c_n$ , so does

$$\sum_{k=1}^n c_k u_k.$$

6. Suppose  $k(\mathbf{x}) = k$ , a constant and  $f = 0$ . Then in one dimension, the heat equation is of the form  $u_t = \alpha u_{xx}$ . Show  $u(x, t) = e^{-\alpha n^2 t} \sin(nx)$  satisfies the heat equation<sup>2</sup>.
7. In a linear, viscous, incompressible fluid, the Cauchy stress is of the form

$$T_{ij}(t, \mathbf{y}) = \lambda \left( \frac{v_{i,j}(t, \mathbf{y}) + v_{j,i}(t, \mathbf{y})}{2} \right)$$

where the comma followed by an index indicates the partial derivative with respect to that variable and  $\mathbf{v}$  is the velocity. Thus

$$v_{i,j} = \frac{\partial v_i}{\partial y_j}$$

---

<sup>2</sup>Fourier, an officer in Napoleon's army studied solutions to the heat equation back in 1813. He was interested in heat flow in cannons. He sought to find solutions by adding up infinitely many solutions of this form. Actually, it was a little more complicated because cannons are not one dimensional but it was the beginning of the study of Fourier series, a topic which fascinated mathematicians for the next 150 years and motivated the development of analysis.

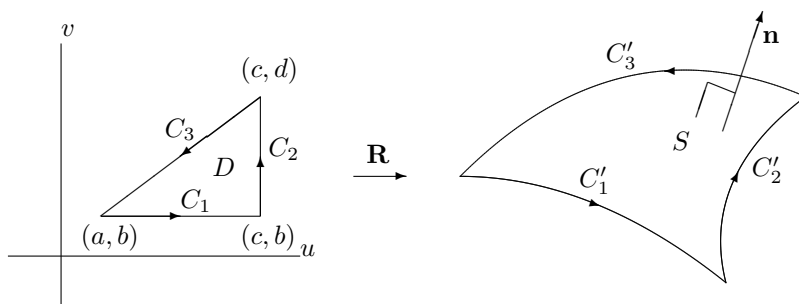
Show, using the balance of mass equation that incompressible implies  $\operatorname{div} \mathbf{v} = 0$ . Next show the balance of momentum equation requires

$$\rho \dot{\mathbf{v}} - \frac{\lambda}{2} \Delta \mathbf{v} = \rho \left[ \frac{\partial \mathbf{v}}{\partial t} + \frac{\partial \mathbf{v}}{\partial y_i} v_i \right] - \frac{\lambda}{2} \Delta \mathbf{v} = \mathbf{b}.$$

This is the famous Navier Stokes equation for incompressible viscous linear fluids. There are still open questions related to this equation, one of which is worth \$1,000,000 at this time.

## 23.7 Stokes Theorem

Stokes theorem relates an integral over the surface to one around the boundary. Consider the following picture.



The triangle on the left is the parameter domain,  $D$ , and the curved surface on the right is the image of this parameter domain,  $S$ , when  $D$  is mapped by  $\mathbf{R}$  to  $\mathbb{R}^3$ . Note the vertices of the triangle are labeled as  $(a,b)$ ,  $(c,b)$ , and  $(c,d)$ . The arrows indicate a direction. Thus moving in the counterclockwise direction around the triangle,  $\mathbf{R}(u,v)$  moves in the indicated direction around the surface,  $S$ . Motion along  $C_1$ , the bottom side of the triangle, corresponds to motion along  $C'_1$  in the indicated direction with similar considerations for the other sides of the parameter domain and  $S$ . The unit vector normal to the parameter domain is coming out of the paper toward you with the motion around the edge of the parameter domain counterclockwise about this unit vector. Think of the fingers of the right hand wrapping in this direction and the thumb pointing in the direction of the unit normal. The unit vector,  $\mathbf{n}$  on the surface,  $S$  has the same relation to motion around the edge of  $S$ , denoted by  $\partial S$ . As in the case of the parameter domain, it is in the direction in which the thumb points when the fingers of the right hand wrap in the indicated direction of motion around  $S$ . Note motion around the parameter domain in the opposite direction results in the unit normal on  $S$  pointing in the opposite direction to the one indicated. It may be useful to imagine the parameter domain drawn on a sheet of rubber lying flat on a table and the surface,  $S$  as what results when one stretches the rubber into a shape in three dimensions. To give a more quantitative description of the unit normal on  $S$ , use the geometric description of the cross product and the observation that  $\mathbf{R}_u(u,v)$  and  $\mathbf{R}_v(u,v)$  are tangent vectors to the surface,  $S$ , at the point  $\mathbf{R}(u,v)$  to conclude that  $\mathbf{n} = \mathbf{R}_u(u,v) \times \mathbf{R}_v(u,v) / |\mathbf{R}_u(u,v) \times \mathbf{R}_v(u,v)|$ . Of course certain assumptions must be made on smoothness. Note first that the normal on  $S$  would not be well defined if  $\mathbf{R}_u(u,v) \times \mathbf{R}_v(u,v) = \mathbf{0}$  and so assume that for all  $(u,v) \in D$ ,  $\mathbf{R}_u(u,v) \times \mathbf{R}_v(u,v) \neq \mathbf{0}$ . It is assumed here that  $\mathbf{R}$  is a  $C^2$  map meaning that all its first and second derivatives exist and are continuous and that  $\mathbf{R}$  is one to one. To minimize confusion in the following argument,  $\mathbf{R}_{,1}$  will denote the partial derivative of  $\mathbf{R}$  with respect to the first variable and  $\mathbf{R}_{,2}$  denotes the partial derivative of  $\mathbf{R}$  with respect to the second variable.

**Lemma 23.7.1** *The area element on  $S$  is of the form,  $dA = |\mathbf{R}_u(u, v) \times \mathbf{R}_v(u, v)| du dv$ .*

**Proof:** This follows from a simple computation. Letting

$$\mathbf{R}(u, v) = (x_1(u, v), x_2(u, v), x_3(u, v)),$$

$$\begin{aligned} |\mathbf{R}_u(u, v) \times \mathbf{R}_v(u, v)|^2 &= \varepsilon_{ijk} x_{j,1} x_{k,2} \varepsilon_{irs} x_{r,1} x_{s,2} \\ &= (\delta_{jr} \delta_{ks} - \delta_{kr} \delta_{js}) x_{j,1} x_{k,2} x_{r,1} x_{s,2} \\ &= x_{j,1} x_{k,2} x_{j,1} x_{k,2} - x_{j,1} x_{k,2} x_{k,1} x_{j,2} \\ &= (\mathbf{R}_u \cdot \mathbf{R}_u)(\mathbf{R}_v \cdot \mathbf{R}_v) - (\mathbf{R}_u \cdot \mathbf{R}_v)^2. \end{aligned}$$

On the other hand,

$$\begin{aligned} D\mathbf{R}^T D\mathbf{R} &= \begin{pmatrix} x_{1,1} & x_{2,1} & x_{3,1} \\ x_{1,2} & x_{2,2} & x_{3,2} \end{pmatrix} \begin{pmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \\ x_{3,1} & x_{3,2} \end{pmatrix} \\ &= \begin{pmatrix} x_{1,1}^2 + x_{2,1}^2 + x_{3,1}^2 & x_{1,1}x_{1,2} + x_{2,1}x_{2,2} + x_{3,1}x_{3,2} \\ x_{1,1}x_{1,2} + x_{2,1}x_{2,2} + x_{3,1}x_{3,2} & x_{1,2}^2 + x_{2,2}^2 + x_{3,2}^2 \end{pmatrix} \end{aligned}$$

and its determinant is  $(\mathbf{R}_u \cdot \mathbf{R}_u)(\mathbf{R}_v \cdot \mathbf{R}_v) - (\mathbf{R}_u \cdot \mathbf{R}_v)^2$  showing that

$$\begin{aligned} dA &= \left( (\mathbf{R}_u \cdot \mathbf{R}_u)(\mathbf{R}_v \cdot \mathbf{R}_v) - (\mathbf{R}_u \cdot \mathbf{R}_v)^2 \right)^{1/2} du dv \\ &= |\mathbf{R}_u(u, v) \times \mathbf{R}_v(u, v)| du dv \end{aligned}$$

as claimed.

**Theorem 23.7.2** *Let  $\mathbf{F}$  be a  $C^1$  vector field defined near  $S$  where the assumptions on  $S$  are as defined above. Then*

$$\int \int_S \nabla \times \mathbf{F} \cdot \mathbf{n} dA = \int_{\partial S} \mathbf{F} \cdot d\mathbf{R} \quad (23.15)$$

**Proof:** The proof is based on the following identity, whose proof is left till later.

$$\begin{aligned} \nabla \times \mathbf{F}(\mathbf{R}(u, v)) \cdot \mathbf{R}_{,1}(u, v) \times \mathbf{R}_{,2}(u, v) &= \\ \frac{\partial \mathbf{F}(\mathbf{R}(u, v))}{\partial u} \cdot \mathbf{R}_{,2}(u, v) - \frac{\partial \mathbf{F}(\mathbf{R}(u, v))}{\partial v} \cdot \mathbf{R}_{,1}(u, v). \end{aligned} \quad (23.16)$$

Using this identity, begin with the left side of (23.15) and arrive at the right side. Note that

$$\mathbf{n} dA = (\mathbf{R}_u(u, v) \times \mathbf{R}_v(u, v)) / |\mathbf{R}_u(u, v) \times \mathbf{R}_v(u, v)| |\mathbf{R}_u(u, v) \times \mathbf{R}_v(u, v)| du dv$$

and so

$$\nabla \times \mathbf{F} \cdot \mathbf{n} dA = \nabla \times \mathbf{F}(\mathbf{R}(u, v)) \cdot \mathbf{R}_{,1}(u, v) \times \mathbf{R}_{,2}(u, v) du dv.$$

Therefore, from (23.16),

$$\begin{aligned} \int \int_S \nabla \times \mathbf{F} \cdot \mathbf{n} dA &= \\ \int \int_D \left( \frac{\partial \mathbf{F}(\mathbf{R}(u, v))}{\partial u} \cdot \mathbf{R}_{,2}(u, v) - \frac{\partial \mathbf{F}(\mathbf{R}(u, v))}{\partial v} \cdot \mathbf{R}_{,1}(u, v) \right) du dv &= \end{aligned}$$

$$\int_b^d \int_{a+(\frac{c-a}{d-b})(v-b)}^c \frac{\partial \mathbf{F}(\mathbf{R}(u, v))}{\partial u} \cdot \mathbf{R}_{,2}(u, v) \, du \, dv \quad (23.17)$$

$$- \int_a^c \int_0^{b+(\frac{d-b}{c-a})(u-a)} \frac{\partial \mathbf{F}(\mathbf{R}(u, v))}{\partial v} \cdot \mathbf{R}_{,1}(u, v) \, dv \, du. \quad (23.18)$$

Integrate the integral of (23.17) by parts. This gives

$$\int_b^d \mathbf{F}(\mathbf{R}(u, v)) \cdot \mathbf{R}_{,2}(u, v) \Big|_{a+(\frac{c-a}{d-b})(v-b)}^c \, dv - \int_b^d \int_{a+(\frac{c-a}{d-b})(v-b)}^c \mathbf{F}(\mathbf{R}(u, v)) \cdot \mathbf{R}_{,21}(u, v) \, du \, dv \quad (23.19)$$

while the expression in (23.18) yields, upon integrating by parts, the following:

$$- \int_a^c \mathbf{F}(\mathbf{R}(u, v)) \cdot \mathbf{R}_{,1}(u, v) \Big|_0^{b+(\frac{d-b}{c-a})(u-a)} + \int_a^c \int_0^{b+(\frac{d-b}{c-a})(u-a)} \mathbf{F}(\mathbf{R}(u, v)) \cdot \mathbf{R}_{,12}(u, v) \, dv \, du. \quad (23.20)$$

The last integrals in (23.19) and (23.20) are the same because of the assumption that  $\mathbf{R}$  is  $C^2$ . Therefore, this reduces to

$$\begin{aligned} & \int_b^d \mathbf{F}(\mathbf{R}(c, v)) \cdot \mathbf{R}_{,2}(c, v) \, dv - \\ & \int_b^d \mathbf{F}\left(\mathbf{R}\left(a + \left(\frac{c-a}{d-b}\right)(v-b), v\right)\right) \cdot \mathbf{R}_{,2}\left(a + \left(\frac{c-a}{d-b}\right)(v-b), v\right) \, dv \\ & - \int_a^c \mathbf{F}\left(\mathbf{R}\left(u, b + \left(\frac{d-b}{c-a}\right)(u-a)\right)\right) \cdot \\ & \mathbf{R}_{,1}\left(u, b + \left(\frac{d-b}{c-a}\right)(u-a)\right) \, du + \int_a^c \mathbf{F}(\mathbf{R}(u, 0)) \cdot \mathbf{R}_{,1}(u, 0) \, du. \end{aligned} \quad (23.21)$$

It remains to recognize this as  $\int_{\partial S} \mathbf{F} \cdot d\mathbf{R}$ . The first integral above equals  $\int_{C'_2} \mathbf{F} \cdot d\mathbf{R}$  and the last integral equals  $\int_{C'_1} \mathbf{F} \cdot d\mathbf{R}$  directly from the definition of these expressions. It only remains to identify the sum of the other two with  $\int_{C'_3} \mathbf{F} \cdot d\mathbf{R}$ . Change the variable in the first of these as  $v = d + t(b-d)$  and in the second by  $u = c + t(a-c)$ . Thus the  $-1$  times the first equals,

$$(b-d) \int_0^1 \mathbf{F}(\mathbf{R}(c+t(a-c), d+t(b-d))) \cdot \mathbf{R}_{,2}(c+t(a-c), d+t(b-d)) \, dt,$$

while  $-1$  times the second equals

$$(a-c) \int_0^1 \mathbf{F}(\mathbf{R}(c+t(a-c), d+t(b-d))) \cdot \mathbf{R}_{,1}(c+t(a-c), d+t(b-d)) \, dt.$$

Now letting  $\tilde{\mathbf{R}}(t) \equiv \mathbf{R}(c+t(a-c), d+t(b-d))$  for  $t \in [0, 1]$ , it follows that  $\tilde{\mathbf{R}}$  is a parameterization for  $C'_3$  and also the chain rule implies

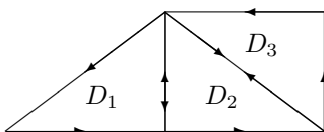
$$\tilde{\mathbf{R}}'(t) = \mathbf{R}_{,1}(c+t(a-c), d+t(b-d))(a-c) + \mathbf{R}_{,2}(c+t(a-c), d+t(b-d))(b-d).$$

Therefore, the sum of the two remaining terms reduces to

$$\int_0^1 \mathbf{F}(\tilde{\mathbf{R}}(t)) \cdot \tilde{\mathbf{R}}'(t) \, dt \equiv \int_{C'_3} \mathbf{F} \cdot d\mathbf{R}$$

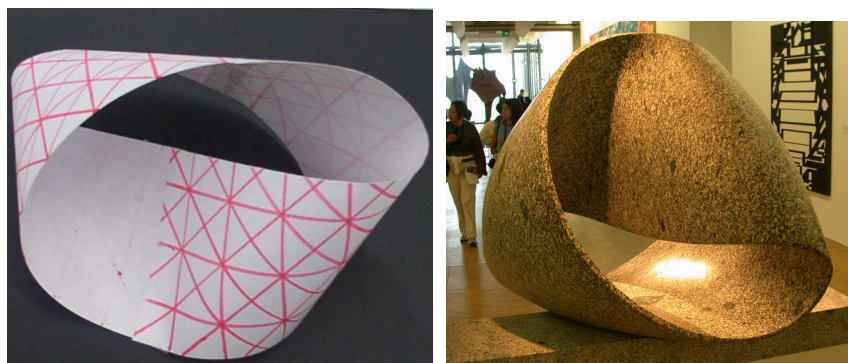
and this shows (23.21) reduces to  $\int_{\partial S} \mathbf{F} \cdot d\mathbf{R}$ . This proves Stoke's theorem.

The above argument was based on a simple parameter domain but other parameter domains work out the same way. In particular, other triangles work out the same. Now consider the following picture of a parameter domain which is made up of simple parameter domains of the sort just considered.



The union of the three triangles shown is considered the parameter domain and  $\mathbf{R}$  maps  $D_i$  to  $S_i$  with the direction of motion around  $\partial S_i$  determined by the direction of motion around  $D_i$  as described above. Then  $\int \int_S \nabla \times \mathbf{F} \cdot \mathbf{n} dA$  equals the sum of the surface integrals over each  $S_i$ . Also,  $\int_{\partial S} \mathbf{F} \cdot d\mathbf{R}$  is equal to the sum of the integrals around  $\partial S_i$  in the indicated direction because the integrals along  $\partial S_i \cap \partial S_j$  cancel due to the fact these are taken in opposite directions. Therefore, knowledge of Stoke's theorem on each piece of  $S$  implies Stoke's theorem on all of  $S$ . The conclusion is that Stoke's theorem holds under very general conditions. As before, think of the parameter domain drawn on a piece of rubber. Stokes theorem will hold for surfaces obtained from the parameter domain by stretching the rubber in three dimensions.

However, you should note that in the above discussion the surface for which Stoke's theorem applies must have two sides. By specifying a direction around the surface, you can let the fingers of your right hand wrap in that direction and the thumb determines the positive side of the surface. Moving around the surface in the other direction yields the other side. Some surfaces have only one side and it is impossible to carry out any treatment of Stoke's theorem for these surfaces. A simple example is the Möbius band, illustrated in the following picture. It is made by putting a twist in a strip of paper and gluing the ends together. Also included is a sculpture depicting the same thing. If you follow the shaded part and keep following it, you will go from the shaded to the unshaded and then back to the shaded and see there is only one side. These Möbius bands have some practical usefulness. In old machine shops where the equipment was run by a belt, they would put such a twist in the belt to spread the surface wear on the belt over twice the area.



Stokes theorem holds under far weaker smoothness assumptions than those presented above. The assumption that the map,  $\mathbf{R}$  is  $C^2$  can be eliminated by approximation of less smooth surfaces with surfaces which are defined in terms of a parameterization,  $\mathbf{R}$  which is  $C^2$ . Using the theory of the Lebesgue integral one can give very general versions of this theorem. The main ideas involve orientation and being able to define surface area.

It remains to verify formula (23.16).

The identity follows from the reduction identity of Lemma 13.6.3 on Page 331 along with the repeated index summation convention discussed there.

$$\begin{aligned}
 \nabla \times \mathbf{F} \cdot \mathbf{R}_{,1} \times \mathbf{R}_{,2} &= \varepsilon_{ijk} (\partial_j F_k) \varepsilon_{irs} R_{r,1} R_{s,2} \\
 &= (\delta_{jr} \delta_{ks} - \delta_{js} \delta_{kr}) (\partial_j F_k) R_{r,1} R_{s,2} \\
 &= (\partial_j F_k) R_{j,1} R_{k,2} - (\partial_j F_k) R_{k,1} R_{j,2} \\
 &= \frac{\partial \mathbf{F}}{\partial u} \cdot \mathbf{R}_{,2} - \frac{\partial \mathbf{F}}{\partial v} \cdot \mathbf{R}_{,1}.
 \end{aligned}$$

Recall the symbol,  $\varepsilon_{ijk}$  is the permutation symbol which is defined by

$$\varepsilon_{ijk} \equiv \begin{cases} 1 & \text{if } (i, j, k) = (1, 2, 3), (2, 3, 1), \text{ or } (3, 1, 2) \\ -1 & \text{if } (i, j, k) = (2, 1, 3), (1, 3, 2), \text{ or } (3, 2, 1) \\ 0 & \text{if } (i, j, k) \text{ has any repeated indices} \end{cases}$$

and  $\delta_{ij}$  is the Kronecker symbol defined by  $\delta_{ij} = 1$  if  $i = j$  and  $\delta_{ij} = 0$  if  $i \neq j$ .

### 23.7.1 Green's Theorem

If the surface,  $S$ , lies in the  $xy$  plane and if the vector field,  $\mathbf{F} = P(x, y)\mathbf{i} + Q(x, y)\mathbf{j}$ , Stokes theorem reduces to Green's theorem. In fact, Stoke's theorem is sometimes proved by using Green's theorem. For this approach, see Vector Analysis by Davis and Snider, [5]. In Green's theorem  $S \subseteq \mathbb{R}^2$  and so the area element is  $dx dy$ . Also, for  $\mathbf{F}$  given as above,  $\nabla \times \mathbf{F} = \left(\frac{\partial Q}{\partial x} - \frac{\partial P}{\partial y}\right)\mathbf{k}$  and the unit normal to  $S$  is just  $\mathbf{k}$  assuming the direction around  $S$  is counter clockwise. Therefore, from Stoke's theorem,

$$\int \int_S \left( \frac{\partial Q}{\partial x} - \frac{\partial P}{\partial y} \right) dx dy = \int \int_S \nabla \times \mathbf{F} \cdot \mathbf{n} dA = \int_{\partial S} \mathbf{F} \cdot d\mathbf{R}.$$



The ends of the above give us Green's theorem. Actually, the line integral on the right is normally written in differential form notation. To do this note that  $d\mathbf{R} = dx\mathbf{i} + dy\mathbf{j}$  and so

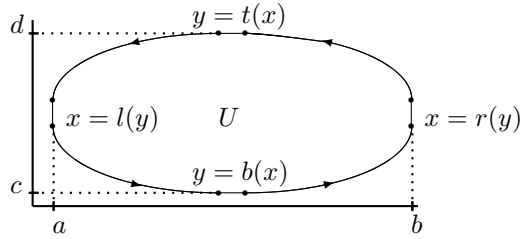
$$\int_{\partial S} \mathbf{F} \cdot d\mathbf{R} = \int_{\partial S} P dx + Q dy$$

but it means the same thing.

## 23.8 Green's Theorem Again

In the above it was shown that Green's theorem follows from Stoke's theorem. You can also go the other way and begin with Green's theorem. I will first establish Green's theorem for regions of a particular sort and then show that the theorem holds for many other regions also. Suppose a region is of the form indicated in the following picture in which

$$\begin{aligned} U &= \{(x, y) : x \in (a, b) \text{ and } y \in (b(x), t(x))\} \\ &= \{(x, y) : y \in (c, d) \text{ and } x \in (l(y), r(y))\}. \end{aligned}$$



I will refer to such a region as being convex in both the  $x$  and  $y$  directions.

**Lemma 23.8.1** Let  $\mathbf{F}(x, y) \equiv (P(x, y), Q(x, y))$  be a  $C^1$  vector field defined near  $U$  where  $U$  is a region of the sort indicated in the above picture which is convex in both the  $x$  and  $y$  directions. Suppose also that the functions,  $r, l, t$ , and  $b$  in the above picture are all  $C^1$  functions and denote by  $\partial U$  the boundary of  $U$  oriented such that the direction of motion is counter clockwise. (As you walk around  $U$  on  $\partial U$ , the points of  $U$  are on your left.) Then

$$\begin{aligned} \int_{\partial U} P dx + Q dy &\equiv \\ \int_{\partial U} \mathbf{F} \cdot d\mathbf{R} &= \iint_U \left( \frac{\partial Q}{\partial x} - \frac{\partial P}{\partial y} \right) dA. \end{aligned} \quad (23.22)$$

**Proof:** First consider the right side of (23.22).

$$\begin{aligned} &\iint_U \left( \frac{\partial Q}{\partial x} - \frac{\partial P}{\partial y} \right) dA \\ &= \int_c^d \int_{l(y)}^{r(y)} \frac{\partial Q}{\partial x} dx dy - \int_a^b \int_{b(x)}^{t(x)} \frac{\partial P}{\partial y} dy dx \\ &= \int_c^d (Q(r(y), y) - Q(l(y), y)) dy + \int_a^b (P(x, b(x)) - P(x, t(x))) dx. \end{aligned} \quad (23.23)$$

Now consider the left side of (23.22). Denote by  $V$  the vertical parts of  $\partial U$  and by  $H$  the horizontal parts.

$$\int_{\partial U} \mathbf{F} \cdot d\mathbf{R} =$$

$$\begin{aligned}
&= \int_{\partial U} (\langle 0, Q \rangle + \langle P, 0 \rangle) \cdot d\mathbf{R} \\
&= \int_c^d \langle 0, Q(r(s), s) \rangle \cdot \langle r'(s), 1 \rangle ds + \int_H \langle 0, Q(r(s), s) \rangle \cdot \langle \pm 1, 0 \rangle ds \\
&\quad - \int_c^d \langle 0, Q(l(s), s) \rangle \cdot \langle l'(s), 1 \rangle ds + \int_a^b \langle P(s, b(s)), 0 \rangle \cdot \langle 1, b'(s) \rangle ds \\
&\quad + \int_V \langle P(s, b(s)), 0 \rangle \cdot \langle 0, \pm 1 \rangle ds - \int_a^b \langle P(s, t(s)), 0 \rangle \cdot \langle 1, t'(s) \rangle ds \\
&= \int_c^d Q(r(s), s) ds - \int_c^d Q(l(s), s) ds + \int_a^b P(s, b(s)) ds - \int_a^b P(s, t(s)) ds
\end{aligned}$$

which coincides with (23.23). This proves the lemma.

**Corollary 23.8.2** *Let everything be the same as in Lemma 23.8.1 but only assume the functions  $r, l, t$ , and  $b$  are continuous and piecewise  $C^1$  functions. Then the conclusion this lemma is still valid.*

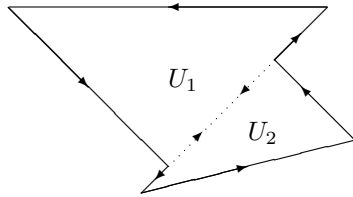
**Proof:** The details are left for you. All you have to do is to break up the various line integrals into the sum of integrals over sub intervals on which the function of interest is  $C^1$ .

From this corollary, it follows (23.22) is valid for any triangle for example.

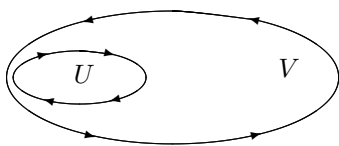
Now suppose (23.22) holds for  $U_1, U_2, \dots, U_m$  and the open sets,  $U_k$  have the property that no two have nonempty intersection and their boundaries intersect only in a finite number of piecewise smooth curves. Then (23.22) must hold for  $U \equiv \cup_{i=1}^m U_i$ , the union of these sets. This is because

$$\begin{aligned}
&\iint_U \left( \frac{\partial Q}{\partial x} - \frac{\partial P}{\partial y} \right) dA = \\
&= \sum_{k=1}^m \iint_{U_k} \left( \frac{\partial Q}{\partial x} - \frac{\partial P}{\partial y} \right) dA \\
&= \sum_{k=1}^m \int_{\partial U_k} \mathbf{F} \cdot d\mathbf{R} = \int_{\partial U} \mathbf{F} \cdot d\mathbf{R}
\end{aligned}$$

because if  $\Gamma = \partial U_k \cap \partial U_j$ , then its orientation as a part of  $\partial U_k$  is opposite to its orientation as a part of  $\partial U_j$  and consequently the line integrals over  $\Gamma$  will cancel, points of  $\Gamma$  also not being in  $\partial U$ . As an illustration, consider the following picture for two such  $U_k$ .



Similarly, if  $U \subseteq V$  and if also  $\partial U \subseteq V$  and both  $U$  and  $V$  are open sets for which (23.22) holds, then the open set,  $V \setminus (U \cup \partial U)$  consisting of what is left in  $V$  after deleting  $U$  along with its boundary also satisfies (23.22). Roughly speaking, you can drill holes in a region for which (23.22) holds and get another region for which this continues to hold provided (23.22) holds for the holes. To see why this is so, consider the following picture which typifies the situation just described.



Then

$$\begin{aligned}\int_{\partial V} \mathbf{F} \cdot d\mathbf{R} &= \int \int_V \left( \frac{\partial Q}{\partial x} - \frac{\partial P}{\partial y} \right) dA \\ &= \int \int_U \left( \frac{\partial Q}{\partial x} - \frac{\partial P}{\partial y} \right) dA + \int \int_{V \setminus U} \left( \frac{\partial Q}{\partial x} - \frac{\partial P}{\partial y} \right) dA \\ &= \int_{\partial U} \mathbf{F} \cdot d\mathbf{R} + \int \int_{V \setminus U} \left( \frac{\partial Q}{\partial x} - \frac{\partial P}{\partial y} \right) dA\end{aligned}$$

and so

$$\int \int_{V \setminus U} \left( \frac{\partial Q}{\partial x} - \frac{\partial P}{\partial y} \right) dA = \int_{\partial V} \mathbf{F} \cdot d\mathbf{R} - \int_{\partial U} \mathbf{F} \cdot d\mathbf{R}$$

which equals

$$\int_{\partial(V \setminus U)} \mathbf{F} \cdot d\mathbf{R}$$

where  $\partial V$  is oriented as shown in the picture. (If you walk around the region,  $V \setminus U$  with the area on the left, you get the indicated orientation for this curve.)

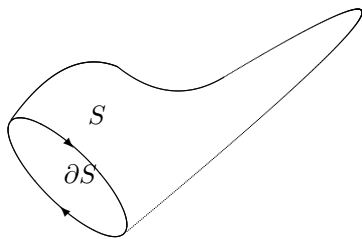
You can see that (23.22) is valid quite generally. This verifies the following theorem.

**Theorem 23.8.3** (*Green's Theorem*) Let  $U$  be an open set in the plane and let  $\partial U$  be piecewise smooth and let  $\mathbf{F}(x, y) = \langle P(x, y), Q(x, y) \rangle$  be a  $C^1$  vector field defined near  $U$ . Then it is often the case that

$$\int_{\partial U} \mathbf{F} \cdot d\mathbf{R} = \int \int_U \left( \frac{\partial Q}{\partial x}(x, y) - \frac{\partial P}{\partial y}(x, y) \right) dA.$$

## 23.9 Stoke's Theorem From Green's Theorem

Stoke's theorem is a generalization of Green's theorem which relates the integral over a surface to the integral around the boundary of the surface. These terms are a little different from what occurs in  $\mathbb{R}^2$ . To describe this, consider a sock. The surface is the sock and its boundary will be the edge of the opening of the sock in which you place your foot. Another way to think of this is to imagine a region in  $\mathbb{R}^2$  of the sort discussed above for Green's theorem. Suppose it is on a sheet of rubber and the sheet of rubber is stretched in three dimensions. The boundary of the resulting surface is the result of the stretching applied to the boundary of the original region in  $\mathbb{R}^2$ . Here is a picture describing the situation.



Recall the following definition of the curl of a vector field.

**Definition 23.9.1** *Let*

$$\mathbf{F}(x, y, z) = \langle F_1(x, y, z), F_2(x, y, z), F_3(x, y, z) \rangle$$

*be a  $C^1$  vector field defined on an open set,  $V$  in  $\mathbb{R}^3$ . Then*

$$\begin{aligned} \nabla \times \mathbf{F} &\equiv \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ \frac{\partial}{\partial x} & \frac{\partial}{\partial y} & \frac{\partial}{\partial z} \\ F_1 & F_2 & F_3 \end{vmatrix} \\ &\equiv \left( \frac{\partial F_3}{\partial y} - \frac{\partial F_2}{\partial z} \right) \mathbf{i} + \left( \frac{\partial F_1}{\partial z} - \frac{\partial F_3}{\partial x} \right) \mathbf{j} + \left( \frac{\partial F_2}{\partial x} - \frac{\partial F_1}{\partial y} \right) \mathbf{k}. \end{aligned}$$

*This is also called curl  $(\mathbf{F})$  and written as indicated,  $\nabla \times \mathbf{F}$ .*

The following lemma gives the fundamental identity which will be used in the proof of Stoke's theorem.

**Lemma 23.9.2** *Let  $\mathbf{R} : U \rightarrow V \subseteq \mathbb{R}^3$  where  $U$  is an open subset of  $\mathbb{R}^2$  and  $V$  is an open subset of  $\mathbb{R}^3$ . Suppose  $\mathbf{R}$  is  $C^2$  and let  $\mathbf{F}$  be a  $C^1$  vector field defined in  $V$ .*

$$(\mathbf{R}_u \times \mathbf{R}_v) \cdot (\nabla \times \mathbf{F})(\mathbf{R}(u, v)) = ((\mathbf{F} \circ \mathbf{R})_u \cdot \mathbf{R}_v - (\mathbf{F} \circ \mathbf{R})_v \cdot \mathbf{R}_u)(u, v). \quad (23.24)$$

**Proof:** Start with the left side and let  $x_i = R_i(u, v)$  for short.

$$\begin{aligned} (\mathbf{R}_u \times \mathbf{R}_v) \cdot (\nabla \times \mathbf{F})(\mathbf{R}(u, v)) &= \varepsilon_{ijk} x_{ju} x_{kv} \varepsilon_{irs} \frac{\partial F_s}{\partial x_r} \\ &= (\delta_{jr} \delta_{ks} - \delta_{js} \delta_{kr}) x_{ju} x_{kv} \frac{\partial F_s}{\partial x_r} \\ &= x_{ju} x_{kv} \frac{\partial F_k}{\partial x_j} - x_{ju} x_{kv} \frac{\partial F_j}{\partial x_k} \\ &= \mathbf{R}_v \cdot \frac{\partial (\mathbf{F} \circ \mathbf{R})}{\partial u} - \mathbf{R}_u \cdot \frac{\partial (\mathbf{F} \circ \mathbf{R})}{\partial v} \end{aligned}$$

which proves (23.24).

The proof of Stoke's theorem given next follows [5]. First, it is convenient to give a definition.

**Definition 23.9.3** *A vector valued function,  $\mathbf{R} : U \subseteq \mathbb{R}^m \rightarrow \mathbb{R}^n$  is said to be in  $C^k(\bar{U}, \mathbb{R}^n)$  if it is the restriction to  $\bar{U}$  of a vector valued function which is defined on  $\mathbb{R}^m$  and is  $C^k$ . That is this function has continuous partial derivatives up to order  $k$ .*

**Theorem 23.9.4 (Stoke's Theorem)** *Let  $U$  be any region in  $\mathbb{R}^2$  for which the conclusion of Green's theorem holds and let  $\mathbf{R} \in C^2(\bar{U}, \mathbb{R}^3)$  be a one to one function satisfying  $|(\mathbf{R}_u \times \mathbf{R}_v)(u, v)| \neq 0$  for all  $(u, v) \in U$  and let  $S$  denote the surface,*

$$\begin{aligned} S &\equiv \{ \mathbf{R}(u, v) : (u, v) \in U \}, \\ \partial S &\equiv \{ \mathbf{R}(u, v) : (u, v) \in \partial U \} \end{aligned}$$

*where the orientation on  $\partial S$  is consistent with the counter clockwise orientation on  $\partial U$  ( $U$  is on the left as you walk around  $\partial U$ ). Then for  $\mathbf{F}$  a  $C^1$  vector field defined near  $S$ ,*

$$\int_{\partial S} \mathbf{F} \cdot d\mathbf{R} = \int \int_S \text{curl}(\mathbf{F}) \cdot \mathbf{n} dS$$

where  $\mathbf{n}$  is the normal to  $S$  defined by

$$\mathbf{n} \equiv \frac{\mathbf{R}_u \times \mathbf{R}_v}{|\mathbf{R}_u \times \mathbf{R}_v|}.$$

**Proof:** Letting  $C$  be an oriented part of  $\partial U$  having parametrization,  $\mathbf{r}(t) \equiv (u(t), v(t))$  for  $t \in [\alpha, \beta]$  and letting  $\mathbf{R}(C)$  denote the oriented part of  $\partial S$  corresponding to  $C$ ,

$$\begin{aligned} \int_{\mathbf{R}(C)} \mathbf{F} \cdot d\mathbf{R} &= \\ &= \int_{\alpha}^{\beta} \mathbf{F}(\mathbf{R}(u(t), v(t))) \cdot (\mathbf{R}_u u'(t) + \mathbf{R}_v v'(t)) dt \\ &= \int_{\alpha}^{\beta} \mathbf{F}(\mathbf{R}(u(t), v(t))) \mathbf{R}_u(u(t), v(t)) u'(t) dt \\ &\quad + \int_{\alpha}^{\beta} \mathbf{F}(\mathbf{R}(u(t), v(t))) \mathbf{R}_v(u(t), v(t)) v'(t) dt \\ &= \int_C \langle (\mathbf{F} \circ \mathbf{R}) \cdot \mathbf{R}_u, (\mathbf{F} \circ \mathbf{R}) \cdot \mathbf{R}_v \rangle \cdot d\mathbf{r}. \end{aligned}$$

Since this holds for each such piece of  $\partial U$ , it follows

$$\int_{\partial S} \mathbf{F} \cdot d\mathbf{R} = \int_{\partial U} \langle (\mathbf{F} \circ \mathbf{R}) \cdot \mathbf{R}_u, (\mathbf{F} \circ \mathbf{R}) \cdot \mathbf{R}_v \rangle \cdot d\mathbf{r}.$$

By the assumption that the conclusion of Green's theorem holds for  $U$ , this equals

$$\begin{aligned} &\int \int_U [((\mathbf{F} \circ \mathbf{R}) \cdot \mathbf{R}_v)_u - ((\mathbf{F} \circ \mathbf{R}) \cdot \mathbf{R}_u)_v] dA \\ &= \int \int_U [(\mathbf{F} \circ \mathbf{R})_u \cdot \mathbf{R}_v + (\mathbf{F} \circ \mathbf{R}) \cdot \mathbf{R}_{vu} - (\mathbf{F} \circ \mathbf{R}) \cdot \mathbf{R}_{uv} - (\mathbf{F} \circ \mathbf{R})_v \cdot \mathbf{R}_u] dA \\ &= \int \int_U [(\mathbf{F} \circ \mathbf{R})_u \cdot \mathbf{R}_v - (\mathbf{F} \circ \mathbf{R})_v \cdot \mathbf{R}_u] dA \end{aligned}$$

the last step holding by equality of mixed partial derivatives, a result of the assumption that  $\mathbf{R}$  is  $C^2$ . Now by Lemma 23.9.2, this equals

$$\begin{aligned} &\int \int_U (\mathbf{R}_u \times \mathbf{R}_v) \cdot (\nabla \times \mathbf{F}) dA \\ &= \int \int_U \nabla \times \mathbf{F} \cdot (\mathbf{R}_u \times \mathbf{R}_v) dA \\ &= \int \int_S \nabla \times \mathbf{F} \cdot \mathbf{n} dS \end{aligned}$$

because  $dS = |(\mathbf{R}_u \times \mathbf{R}_v)| dA$  and  $\mathbf{n} = \frac{(\mathbf{R}_u \times \mathbf{R}_v)}{|(\mathbf{R}_u \times \mathbf{R}_v)|}$ . Thus

$$\begin{aligned} (\mathbf{R}_u \times \mathbf{R}_v) dA &= \frac{(\mathbf{R}_u \times \mathbf{R}_v)}{|(\mathbf{R}_u \times \mathbf{R}_v)|} |(\mathbf{R}_u \times \mathbf{R}_v)| dA \\ &= \mathbf{n} dS. \end{aligned}$$

This proves Stoke's theorem.

### 23.9.1 Conservative Vector Fields

**Definition 23.9.5** A vector field,  $\mathbf{F}$  defined in a three dimensional region is said to be conservative<sup>3</sup> if for every piecewise smooth closed curve,  $C$ , it follows  $\int_C \mathbf{F} \cdot d\mathbf{R} = 0$ .

**Definition 23.9.6** Let  $(\mathbf{x}, \mathbf{p}_1, \dots, \mathbf{p}_n, \mathbf{y})$  be an ordered list of points in  $\mathbb{R}^p$ . Let

$$\mathbf{p}(\mathbf{x}, \mathbf{p}_1, \dots, \mathbf{p}_n, \mathbf{y})$$

denote the piecewise smooth curve consisting of a straight line segment from  $\mathbf{x}$  to  $\mathbf{p}_1$  and then the straight line segment from  $\mathbf{p}_1$  to  $\mathbf{p}_2 \dots$  and finally the straight line segment from  $\mathbf{p}_n$  to  $\mathbf{y}$ . This is called a polygonal curve. An open set in  $\mathbb{R}^p$ ,  $U$ , is said to be a region if it has the property that for any two points,  $\mathbf{x}, \mathbf{y} \in U$ , there exists a polygonal curve joining the two points.

Conservative vector fields are important because of the following theorem.

**Theorem 23.9.7** Let  $U$  be a region in  $\mathbb{R}^p$  and let  $\mathbf{F} : U \rightarrow \mathbb{R}^p$  be a continuous vector field. Then  $\mathbf{F}$  is conservative if and only if there exists a scalar valued function of  $p$  variables,  $\phi$  such that  $\mathbf{F} = \nabla \phi$ . Furthermore, if  $C$  is an oriented curve which goes from  $\mathbf{x}$  to  $\mathbf{y}$  in  $U$ , then

$$\int_C \mathbf{F} \cdot d\mathbf{R} = \phi(\mathbf{y}) - \phi(\mathbf{x}). \quad (23.25)$$

Thus the line integral is path independent in this case. This function,  $\phi$  is called a scalar potential for  $\mathbf{F}$ .

**Proof:** To save space and fussing over things which are unimportant, denote by  $\mathbf{p}(\mathbf{x}_0, \mathbf{x})$  a polygonal curve from  $\mathbf{x}_0$  to  $\mathbf{x}$ . Thus the orientation is such that it goes from  $\mathbf{x}_0$  to  $\mathbf{x}$ . The curve  $\mathbf{q}(\mathbf{x}, \mathbf{x}_0)$  denotes the same set of points but in the opposite order. Suppose first  $\mathbf{F}$  is conservative. Fix  $\mathbf{x}_0 \in U$  and let

$$\phi(\mathbf{x}) \equiv \int_{\mathbf{p}(\mathbf{x}_0, \mathbf{x})} \mathbf{F} \cdot d\mathbf{R}.$$

This is well defined because if  $\mathbf{q}(\mathbf{x}_0, \mathbf{x})$  is another polygonal curve joining  $\mathbf{x}_0$  to  $\mathbf{x}$ , Then the curve obtained by following  $\mathbf{p}(\mathbf{x}_0, \mathbf{x})$  from  $\mathbf{x}_0$  to  $\mathbf{x}$  and then from  $\mathbf{x}$  to  $\mathbf{x}_0$  along  $\mathbf{q}(\mathbf{x}, \mathbf{x}_0)$  is a closed piecewise smooth curve and so by assumption, the line integral along this closed curve equals 0. However, this integral is just

$$\int_{\mathbf{p}(\mathbf{x}_0, \mathbf{x})} \mathbf{F} \cdot d\mathbf{R} + \int_{\mathbf{q}(\mathbf{x}, \mathbf{x}_0)} \mathbf{F} \cdot d\mathbf{R} = \int_{\mathbf{p}(\mathbf{x}_0, \mathbf{x})} \mathbf{F} \cdot d\mathbf{R} - \int_{\mathbf{q}(\mathbf{x}_0, \mathbf{x})} \mathbf{F} \cdot d\mathbf{R}$$

which shows

$$\int_{\mathbf{p}(\mathbf{x}_0, \mathbf{x})} \mathbf{F} \cdot d\mathbf{R} = \int_{\mathbf{q}(\mathbf{x}_0, \mathbf{x})} \mathbf{F} \cdot d\mathbf{R}$$

and that  $\phi$  is well defined. For small  $t$ ,

$$\begin{aligned} \frac{\phi(\mathbf{x} + t\mathbf{e}_i) - \phi(\mathbf{x})}{t} &= \frac{\int_{\mathbf{p}(\mathbf{x}_0, \mathbf{x} + t\mathbf{e}_i)} \mathbf{F} \cdot d\mathbf{R} - \int_{\mathbf{p}(\mathbf{x}_0, \mathbf{x})} \mathbf{F} \cdot d\mathbf{R}}{t} \\ &= \frac{\int_{\mathbf{p}(\mathbf{x}_0, \mathbf{x})} \mathbf{F} \cdot d\mathbf{R} + \int_{\mathbf{p}(\mathbf{x}, \mathbf{x} + t\mathbf{e}_i)} \mathbf{F} \cdot d\mathbf{R} - \int_{\mathbf{p}(\mathbf{x}_0, \mathbf{x})} \mathbf{F} \cdot d\mathbf{R}}{t}. \end{aligned}$$

---

<sup>3</sup>There is no such thing as a liberal vector field.

Since  $U$  is open, for small  $t$ , the ball of radius  $|t|$  centered at  $\mathbf{x}$  is contained in  $U$ . Therefore, the line segment from  $\mathbf{x}$  to  $\mathbf{x} + t\mathbf{e}_i$  is also contained in  $U$  and so one can take  $\mathbf{p}(\mathbf{x}, \mathbf{x} + t\mathbf{e}_i)(s) = \mathbf{x} + s(t\mathbf{e}_i)$  for  $s \in [0, 1]$ . Therefore, the above difference quotient reduces to

$$\begin{aligned} \frac{1}{t} \int_0^1 \mathbf{F}(\mathbf{x} + s(t\mathbf{e}_i)) \cdot t\mathbf{e}_i \, ds &= \int_0^1 F_i(\mathbf{x} + s(t\mathbf{e}_i)) \, ds \\ &= F_i(\mathbf{x} + s_t(t\mathbf{e}_i)) \end{aligned}$$

by the mean value theorem for integrals. Here  $s_t$  is some number between 0 and 1. By continuity of  $\mathbf{F}$ , this converges to  $F_i(\mathbf{x})$  as  $t \rightarrow 0$ . Therefore,  $\nabla\phi = \mathbf{F}$  as claimed.

Conversely, if  $\nabla\phi = \mathbf{F}$ , then if  $\mathbf{R} : [a, b] \rightarrow \mathbb{R}^p$  is any  $C^1$  curve joining  $\mathbf{x}$  to  $\mathbf{y}$ ,

$$\begin{aligned} \int_a^b \mathbf{F}(\mathbf{R}(t)) \cdot \mathbf{R}'(t) \, dt &= \int_a^b \nabla\phi(\mathbf{R}(t)) \cdot \mathbf{R}'(t) \, dt \\ &= \int_a^b \frac{d}{dt}(\phi(\mathbf{R}(t))) \, dt \\ &= \phi(\mathbf{R}(b)) - \phi(\mathbf{R}(a)) \\ &= \phi(\mathbf{y}) - \phi(\mathbf{x}) \end{aligned}$$

and this verifies (23.25) in the case where the curve joining the two points is smooth. The general case follows immediately from this by using this result on each of the pieces of the piecewise smooth curve. For example if the curve goes from  $\mathbf{x}$  to  $\mathbf{p}$  and then from  $\mathbf{p}$  to  $\mathbf{y}$ , the above would imply the integral over the curve from  $\mathbf{x}$  to  $\mathbf{p}$  is  $\phi(\mathbf{p}) - \phi(\mathbf{x})$  while from  $\mathbf{p}$  to  $\mathbf{y}$  the integral would yield  $\phi(\mathbf{y}) - \phi(\mathbf{p})$ . Adding these gives  $\phi(\mathbf{y}) - \phi(\mathbf{x})$ . The formula (23.25) implies the line integral over any closed curve equals zero because the starting and ending points of such a curve are the same. This proves the theorem.

**Example 23.9.8** Let  $\mathbf{F}(x, y, z) = (\cos x - yz \sin(xz), \cos(xz), -yx \sin(xz))$ . Let  $C$  be a piecewise smooth curve which goes from  $(\pi, 1, 1)$  to  $(\frac{\pi}{2}, 3, 2)$ . Find  $\int_C \mathbf{F} \cdot d\mathbf{R}$ .

The specifics of the curve are not given so the problem is nonsense unless the vector field is conservative. Therefore, it is reasonable to look for the function,  $\phi$  satisfying  $\nabla\phi = \mathbf{F}$ . Such a function satisfies

$$\phi_x = \cos x - y(\sin xz)z$$

and so, assuming  $\phi$  exists,

$$\phi(x, y, z) = \sin x + y \cos(xz) + \psi(y, z).$$

I have to add in the most general thing possible,  $\psi(y, z)$  to ensure possible solutions are not being thrown out. It wouldn't be good at this point to add in a constant since the answer could involve a function of either or both of the other variables. Now from what was just obtained,

$$\phi_y = \cos(xz) + \psi_y = \cos xz$$

and so it is possible to take  $\psi_y = 0$ . Consequently,  $\phi$ , if it exists is of the form

$$\phi(x, y, z) = \sin x + y \cos(xz) + \psi(z).$$

Now differentiating this with respect to  $z$  gives

$$\phi_z = -yx \sin(xz) + \psi_z = -yx \sin(xz)$$

and this shows  $\psi$  does not depend on  $z$  either. Therefore, it suffices to take  $\psi = 0$  and

$$\phi(x, y, z) = \sin(x) + y \cos(xz).$$

Therefore, the desired line integral equals

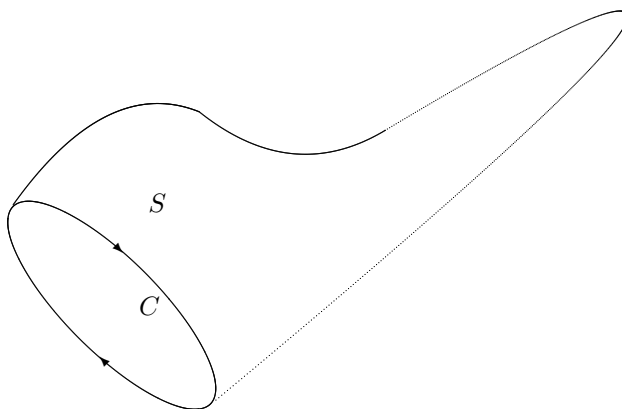
$$\sin\left(\frac{\pi}{2}\right) + 3 \cos(\pi) - (\sin(\pi) + \cos(\pi)) = -1.$$

The above process for finding  $\phi$  will not lead you astray in the case where there does not exist a scalar potential. As an example, consider the following.

**Example 23.9.9** Let  $\mathbf{F}(x, y, z) = (x, y^2x, z)$ . Find a scalar potential for  $\mathbf{F}$  if it exists.

If  $\phi$  exists, then  $\phi_x = x$  and so  $\phi = \frac{x^2}{2} + \psi(y, z)$ . Then  $\phi_y = \psi_y(y, z) = xy^2$  but this is impossible because the left side depends only on  $y$  and  $z$  while the right side depends also on  $x$ . Therefore, this vector field is not conservative and there does not exist a scalar potential.

**Definition 23.9.10** A set of points in three dimensional space,  $V$  is simply connected if every piecewise smooth closed curve,  $C$  is the edge of a surface,  $S$  which is contained entirely within  $V$  in such a way that Stokes theorem holds for the surface,  $S$  and its edge,  $C$ .



This is like a sock. The surface is the sock and the curve,  $C$  goes around the opening of the sock.

As an application of Stoke's theorem, here is a useful theorem which gives a way to check whether a vector field is conservative.

**Theorem 23.9.11** For a three dimensional simply connected open set,  $V$  and  $\mathbf{F}$  a  $C^1$  vector field defined in  $V$ ,  $\mathbf{F}$  is conservative if  $\nabla \times \mathbf{F} = \mathbf{0}$  in  $V$ .

**Proof:** If  $\nabla \times \mathbf{F} = \mathbf{0}$  then taking an arbitrary closed curve,  $C$ , and letting  $S$  be a surface bounded by  $C$  which is contained in  $V$ , Stoke's theorem implies

$$0 = \int_S \nabla \times \mathbf{F} \cdot \mathbf{n} dA = \int_C \mathbf{F} \cdot d\mathbf{R}.$$

Thus  $\mathbf{F}$  is conservative.



### 23.9.2 Maxwell's Equations And The Wave Equation

Many of the ideas presented above are useful in analyzing Maxwell's equations. These equations are derived in advanced physics courses. They are

$$\nabla \times \mathbf{E} + \frac{1}{c} \frac{\partial \mathbf{B}}{\partial t} = \mathbf{0} \quad (23.26)$$

$$\nabla \cdot \mathbf{E} = 4\pi\rho \quad (23.27)$$

$$\nabla \times \mathbf{B} - \frac{1}{c} \frac{\partial \mathbf{E}}{\partial t} = \frac{4\pi}{c} \mathbf{f} \quad (23.28)$$

$$\nabla \cdot \mathbf{B} = 0 \quad (23.29)$$

and it is assumed these hold on all of  $\mathbb{R}^3$  to eliminate technical considerations having to do with whether something is simply connected.

In these equations,  $\mathbf{E}$  is the electrostatic field and  $\mathbf{B}$  is the magnetic field while  $\rho$  and  $\mathbf{f}$  are sources. By (23.29)  $\mathbf{B}$  has a vector potential,  $\mathbf{A}_1$  such that  $\mathbf{B} = \nabla \times \mathbf{A}_1$ . Now go to (23.26) and write

$$\nabla \times \mathbf{E} + \frac{1}{c} \nabla \times \frac{\partial \mathbf{A}_1}{\partial t} = \mathbf{0}$$

showing that

$$\nabla \times \left( \mathbf{E} + \frac{1}{c} \frac{\partial \mathbf{A}_1}{\partial t} \right) = \mathbf{0}$$

It follows  $\mathbf{E} + \frac{1}{c} \frac{\partial \mathbf{A}_1}{\partial t}$  has a scalar potential,  $\psi_1$  satisfying

$$\nabla \psi_1 = \mathbf{E} + \frac{1}{c} \frac{\partial \mathbf{A}_1}{\partial t}. \quad (23.30)$$

Now suppose  $\phi$  is a time dependent scalar field satisfying

$$\nabla^2 \phi - \frac{1}{c^2} \frac{\partial^2 \phi}{\partial t^2} = \frac{1}{c} \frac{\partial \psi_1}{\partial t} - \nabla \cdot \mathbf{A}_1. \quad (23.31)$$

Next define

$$\mathbf{A} \equiv \mathbf{A}_1 + \nabla \phi, \quad \psi \equiv \psi_1 + \frac{1}{c} \frac{\partial \phi}{\partial t}. \quad (23.32)$$

Therefore, in terms of the new variables, (23.31) becomes

$$\nabla^2 \phi - \frac{1}{c^2} \frac{\partial^2 \phi}{\partial t^2} = \frac{1}{c} \left( \frac{\partial \psi}{\partial t} - \frac{1}{c} \frac{\partial^2 \phi}{\partial t^2} \right) - \nabla \cdot \mathbf{A} + \nabla^2 \phi$$

which yields

$$0 = \frac{\partial \psi}{\partial t} - c \nabla \cdot \mathbf{A}. \quad (23.33)$$

Then it follows from Theorem 23.1.3 on Page 594 that  $\mathbf{A}$  is also a vector potential for  $\mathbf{B}$ . That is

$$\nabla \times \mathbf{A} = \mathbf{B}. \quad (23.34)$$

From (23.30)

$$\nabla \left( \psi - \frac{1}{c} \frac{\partial \phi}{\partial t} \right) = \mathbf{E} + \frac{1}{c} \left( \frac{\partial \mathbf{A}}{\partial t} - \nabla \frac{\partial \phi}{\partial t} \right)$$

and so

$$\nabla \psi = \mathbf{E} + \frac{1}{c} \frac{\partial \mathbf{A}}{\partial t}. \quad (23.35)$$

Using (23.28) and (23.35),

$$\nabla \times (\nabla \times \mathbf{A}) - \frac{1}{c} \frac{\partial}{\partial t} \left( \nabla \psi - \frac{1}{c} \frac{\partial \mathbf{A}}{\partial t} \right) = \frac{4\pi}{c} \mathbf{f}. \quad (23.36)$$

Now from Theorem 23.1.3 on Page 594 this implies

$$\nabla (\nabla \cdot \mathbf{A}) - \nabla^2 \mathbf{A} - \nabla \left( \frac{1}{c} \frac{\partial \psi}{\partial t} \right) + \frac{1}{c^2} \frac{\partial^2 \mathbf{A}}{\partial t^2} = \frac{4\pi}{c} \mathbf{f}$$

and using (23.33), this gives

$$\frac{1}{c^2} \frac{\partial^2 \mathbf{A}}{\partial t^2} - \nabla^2 \mathbf{A} = \frac{4\pi}{c} \mathbf{f}. \quad (23.37)$$

Also from (23.35), (23.27), and (23.33),

$$\begin{aligned} \nabla^2 \psi &= \nabla \cdot \mathbf{E} + \frac{1}{c} \frac{\partial}{\partial t} (\nabla \cdot \mathbf{A}) \\ &= 4\pi\rho + \frac{1}{c^2} \frac{\partial^2 \psi}{\partial t^2} \end{aligned}$$

and so

$$\frac{1}{c^2} \frac{\partial^2 \psi}{\partial t^2} - \nabla^2 \psi = -4\pi\rho. \quad (23.38)$$

This is very interesting. If a solution to the wave equations, (23.38), and (23.37) can be found along with a solution to (23.33), then letting the magnetic field be given by (23.34) and letting  $\mathbf{E}$  be given by (23.35) the result is a solution to Maxwells equations. This is significant because wave equations are easier to think of than Maxwell's equations. Note the above argument also showed that it is always possible, by solving another wave equation, to get (23.33) to hold.

## 23.10 Exercises

- Determine whether the vector field,  $(2xy^3 \sin z^4, 3x^2y^2 \sin z^4 + 1, 4x^2y^3 (\cos z^4) z^3 + 1)$  is conservative. If it is conservative, find a potential function.
- Determine whether the vector field,  $(2xy^3 \sin z + y^2 + z, 3x^2y^2 \sin z + 2xy, x^2y^3 \cos z + x)$  is conservative. If it is conservative, find a potential function.
- Determine whether the vector field,  $(2xy^3 \sin z + z, 3x^2y^2 \sin z + 2xy, x^2y^3 \cos z + x)$  is conservative. If it is conservative, find a potential function.
- Find scalar potentials for the following vector fields if it is possible to do so. If it is not possible to do so, explain why.
  - $(y^2, 2xy + \sin z, 2z + y \cos z)$
  - $(2z (\cos (x^2 + y^2)) x, 2z (\cos (x^2 + y^2)) y, \sin (x^2 + y^2) + 2z)$
  - $(f(x), g(y), h(z))$
  - $(xy, z^2, y^3)$
  - $\left( z + 2 \frac{x}{x^2 + y^2 + 1}, 2 \frac{y}{x^2 + y^2 + 1}, x + 3z^2 \right)$

5. If a vector field is not conservative on the set  $U$ , is it possible the same vector field could be conservative on some subset of  $U$ ? Explain and give examples if it is possible. If it is not possible also explain why.
6. Prove that if a vector field,  $\mathbf{F}$  has a scalar potential, then it has infinitely many scalar potentials.
7. Here is a vector field:  $\mathbf{F} \equiv (2xy, x^2 - 5y^4, 3z^2)$ . Find  $\int_C \mathbf{F} \cdot d\mathbf{R}$  where  $C$  is a curve which goes from  $(1, 2, 3)$  to  $(4, -2, 1)$ .
8. Here is a vector field:  $\mathbf{F} \equiv (2xy, x^2 - 5y^4, 3(\cos z^3)z^2)$ . Find  $\int_C \mathbf{F} \cdot d\mathbf{R}$  where  $C$  is a curve which goes from  $(1, 0, 1)$  to  $(-4, -2, 1)$ .
9. Find  $\int_{\partial U} \mathbf{F} \cdot d\mathbf{R}$  where  $U$  is the set,  $\{(x, y) : 2 \leq x \leq 4, 0 \leq y \leq x\}$  and  $\mathbf{F}(x, y) = (x \sin y, y \sin x)$ .
10. Find  $\int_{\partial U} \mathbf{F} \cdot d\mathbf{R}$  where  $U$  is the set,  $\{(x, y) : 2 \leq x \leq 3, 0 \leq y \leq x^2\}$  and  $\mathbf{F}(x, y) = (x \cos y, y + x)$ .
11. Find  $\int_{\partial U} \mathbf{F} \cdot d\mathbf{R}$  where  $U$  is the set,  $\{(x, y) : 1 \leq x \leq 2, x \leq y \leq 3\}$  and  $\mathbf{F}(x, y) = (x \sin y, y \sin x)$ .
12. Find  $\int_{\partial U} \mathbf{F} \cdot d\mathbf{R}$  where  $U$  is the set,  $\{(x, y) : x^2 + y^2 \leq 2\}$  and  $\mathbf{F}(x, y) = (-y^3, x^3)$ .
13. Show that for many open sets in  $\mathbb{R}^2$ , Area of  $U = \int_{\partial U} x dy$ , and Area of  $U = \int_{\partial U} -y dx$  and Area of  $U = \frac{1}{2} \int_{\partial U} -y dx + x dy$ . **Hint:** Use Green's theorem.
14. Let  $P$  be a polygon with vertices  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n), (x_1, y_1)$  encountered as you move over the boundary of the polygon in the counter clockwise direction. Using Problem 13, find a nice formula for the area of the polygon in terms of the vertices.
15. Two smooth oriented surfaces,  $S_1$  and  $S_2$  intersect in a piecewise smooth oriented closed curve,  $C$ . Let  $\mathbf{F}$  be a  $C^1$  vector field defined on  $\mathbb{R}^3$ . Explain why  $\int_{S_1} \text{curl}(\mathbf{F}) \cdot \mathbf{n} dS = \int_{S_2} \text{curl}(\mathbf{F}) \cdot \mathbf{n} dS$ . Here  $\mathbf{n}$  is the normal to the surface which corresponds to the given orientation of the curve,  $C$ .
16. Parametric equations for one arch of a cycloid are given by  $x = a(t - \sin t)$  and  $y = a(1 - \cos t)$  where here  $t \in [0, 2\pi]$ . Sketch a rough graph of this arch of a cycloid and then find the area between this arch and the  $x$  axis. **Hint:** This is very easy using Green's theorem and the vector field,  $\mathbf{F} = (-y, x)$ .
17. Consider the vector field,  $\left(\frac{-y}{(x^2+y^2)}, \frac{x}{(x^2+y^2)}, 0\right) = \mathbf{F}$ . Show that  $\nabla \times \mathbf{F} = \mathbf{0}$  but that for the closed curve, whose parameterization is  $\mathbf{R}(t) = (\cos t, \sin t, 0)$  for  $t \in [0, 2\pi]$ ,  $\int_C \mathbf{F} \cdot d\mathbf{R} \neq 0$ . Therefore, the vector field is not conservative. Does this contradict Theorem 23.9.11? Explain.
18. The cylinder  $x^2 + y^2 = 4$  is intersected with the plane  $x + y + z = 2$ . This yields a closed curve,  $C$ . Orient this curve in the counter clockwise direction when viewed from a point high on the  $z$  axis. Let  $\mathbf{F} = (x^2 y, z + y, x^2)$ . Find  $\int_C \mathbf{F} \cdot d\mathbf{R}$ .
19. The cylinder  $x^2 + 4y^2 = 4$  is intersected with the plane  $x + 3y + 2z = 1$ . This yields a closed curve,  $C$ . Orient this curve in the counter clockwise direction when viewed from a point high on the  $z$  axis. Let  $\mathbf{F} = (y, z + y, x^2)$ . Find  $\int_C \mathbf{F} \cdot d\mathbf{R}$ .

20. The cylinder  $x^2 + y^2 = 4$  is intersected with the plane  $x + 3y + 2z = 1$ . This yields a closed curve,  $C$ . Orient this curve in the clockwise direction when viewed from a point high on the  $z$  axis. Let  $\mathbf{F} = (y, z + y, x)$ . Find  $\int_C \mathbf{F} \cdot d\mathbf{R}$ .
21. Let  $\mathbf{F} = (xz, z^2(y + \sin x), z^3y)$ . Find the surface integral,  $\int \int_S \text{curl}(\mathbf{F}) \cdot \mathbf{n} dA$  where  $S$  is the surface,  $z = 4 - (x^2 + y^2)$ ,  $z \geq 0$ .
22. Let  $\mathbf{F} = (xz, (y^3 + x), z^3y)$ . Find the surface integral,  $\int \int_S \text{curl}(\mathbf{F}) \cdot \mathbf{n} dA$  where  $S$  is the surface,  $z = 16 - (x^2 + y^2)$ ,  $z \geq 0$ .
23. The cylinder  $z = y^2$  intersects the surface  $z = 8 - x^2 - 4y^2$  in a curve,  $C$  which is oriented in the counter clockwise direction when viewed high on the  $z$  axis. Find  $\int_C \mathbf{F} \cdot d\mathbf{R}$  if  $\mathbf{F} = \left(\frac{z^2}{2}, xy, xz\right)$ . **Hint:** This is not too hard if you show you can use Stokes theorem on a domain in the  $xy$  plane.
24. Suppose solutions have been found to (23.38), (23.37), and (23.33). Then define  $\mathbf{E}$  and  $\mathbf{B}$  using (23.35) and (23.34). Verify Maxwell's equations hold for  $\mathbf{E}$  and  $\mathbf{B}$ .
25. Suppose now you have found solutions to (23.38) and (23.37),  $\psi_1$  and  $A_1$ . Then go show again that if  $\phi$  satisfies (23.31) and  $\psi \equiv \psi_1 + \frac{1}{c} \frac{\partial \phi}{\partial t}$ , while  $\mathbf{A} \equiv \mathbf{A}_1 + \nabla \phi$ , then (23.33) holds for  $\mathbf{A}$  and  $\psi$ .
26. Why consider Maxwell's equations? Why not just consider (23.38), (23.37), and (23.33)?
27. Tell which open sets are simply connected.
  - (a) The inside of a car radiator.
  - (b) A donut.
  - (c) The solid part of a cannon ball which contains a void on the interior.
  - (d) The inside of a donut which has had a large bite taken out of it.
  - (e) All of  $\mathbb{R}^3$  except the  $z$  axis.
  - (f) All of  $\mathbb{R}^3$  except the  $xy$  plane.

# The Fundamental Theorem Of Algebra

The fundamental theorem of algebra states that every non constant polynomial having coefficients in  $\mathbb{C}$  has a zero in  $\mathbb{C}$ . If  $\mathbb{C}$  is replaced by  $\mathbb{R}$ , this is not true because of the example,  $x^2 + 1 = 0$ . This theorem is a very remarkable result and notwithstanding its title, all the best proofs depend on either analysis or topology. It was first proved by Gauss in 1797. The proof given here follows Rudin [14]. See also Hardy [9] for a similar proof, more discussion and references. The best proof of this fundamental theorem is based on Liouville's theorem and is found in nearly every book on complex analysis.

**Lemma A.0.1** *Let  $a_k \in \mathbb{C}$  for  $k = 1, \dots, n$  and let  $p(z) \equiv \sum_{k=1}^n a_k z^k$ . Then  $p$  is continuous.*

**Proof:**

$$|az^n - aw^n| \leq |a| |z - w| |z^{n-1} + z^{n-2}w + \dots + w^{n-1}|.$$

Then for  $|z - w| < 1$ , the triangle inequality implies  $|w| < 1 + |z|$  and so if  $|z - w| < 1$ ,

$$|az^n - aw^n| \leq |a| |z - w| n (1 + |z|)^n.$$

If  $\varepsilon > 0$  is given, let

$$\delta < \min \left( 1, \frac{\varepsilon}{|a| n (1 + |z|)^n} \right).$$

It follows from the above inequality that for  $|z - w| < \delta$ ,  $|az^n - aw^n| < \varepsilon$ . The function of the lemma is just the sum of functions of this sort and so it follows that it is also continuous.

**Theorem A.0.2** (*Fundamental theorem of Algebra*) *Let  $p(z)$  be a non constant polynomial. Then there exists  $z \in \mathbb{C}$  such that  $p(z) = 0$ .*

**Proof:** Suppose not. Then

$$p(z) = \sum_{k=0}^n a_k z^k$$

where  $a_n \neq 0$ ,  $n > 0$ . Then

$$|p(z)| \geq |a_n| |z|^n - \sum_{k=0}^{n-1} |a_k| |z|^k$$

and so

$$\lim_{|z| \rightarrow \infty} |p(z)| = \infty. \quad (1.1)$$

Now let

$$\lambda \equiv \inf \{|p(z)| : z \in \mathbb{C}\}.$$

By (1.1), there exists an  $R > 0$  such that if  $|z| > R$ , it follows that  $|p(z)| > \lambda + 1$ . Therefore,

$$\lambda \equiv \inf \{|p(z)| : z \in \mathbb{C}\} = \inf \{|p(z)| : |z| \leq R\}.$$

The set  $\{z : |z| \leq R\}$  is a closed and bounded set in  $\mathbb{R}^2$  and so this infimum is achieved at some point  $w$  with  $|w| \leq R$ . If  $|p(w)| = 0$ , this is a contradiction so assume  $|p(w)| > 0$ . Then consider

$$q(z) \equiv \frac{p(z+w)}{p(w)}.$$

Since  $q(0) = 1$ , it follows  $q(z)$  is of the form

$$q(z) = 1 + c_k z^k + \cdots + c_n z^n$$

where  $c_k$  is the first coefficient which is nonzero. It is also true that  $|q(z)| \geq 1$  by the assumption that  $|p(w)|$  is the smallest value of  $|p(z)|$ .

Therefore,

$$q(z) = 1 + z^k (c_k + d(z))$$

where  $\lim_{z \rightarrow 0} d(z) = 0$ . Now by De Moivre's theorem, Theorem 4.0.4 on Page 74, there exists  $z_n$  such that

$$z_n^k = \frac{\overline{c_k}}{|c_k|^2} \left( \frac{-1}{n} \right).$$

Then by De Moivre's theorem again,  $\lim_{n \rightarrow \infty} d(z_n) = 0$  and

$$q(z_n) = 1 - \frac{1}{n} + \left( \frac{\overline{c_k}}{|c_k|^2} \left( \frac{-1}{n} \right) d(z_n) \right).$$

Therefore, for  $n$  large enough

$$\left| \frac{\overline{c_k}}{|c_k|^2} d(z_n) \right| < \frac{1}{2}$$

and so for such  $n$ ,

$$|q(z_n)| \leq 1 - \frac{1}{n} + \frac{1}{2n} < 1,$$

a contradiction to  $|q(z)| \geq 1$ .

# Bibliography

- [1] **Apostol, T. M.**, *Mathematical Analysis*, Addison Wesley Publishing Co., 1974.
- [2] **Apostol, T. M.**, *Calculus second edition*, Wiley, 1967.
- [3] **Bartle R.G.**, *A Modern Theory of Integration*, Grad. Studies in Math., Amer. Math. Society, Providence, RI, 2000.
- [4] **Chahal J. S.** , *Historical Perspective of Mathematics* 2000 B.C. - 2000 A.D.
- [5] **Davis H. and Snider A.**, *Vector Analysis* Wm. C. Brown 1995.
- [6] **Euclid**, *The Thirteen Books of the Elements*, Dover, 1956.
- [7] **Fitzpatrick P. M.**, *Advanced Calculus a course in Mathematical Analysis*, PWS Publishing Company 1996.
- [8] **Fleming W.**, *Functions of Several Variables*, Springer Verlag 1976.
- [9] **Hardy G.**, *A Course Of Pure Mathematics, Tenth edition*, Cambridge University Press 1992.
- [10] **Kuttler K. L.**, *Basic Analysis*, Rinton
- [11] **Kuttler K.L.**, *Modern Analysis* CRC Press 1998.
- [12] **Lang S.** *Real and Functional analysis* third edition Springer Verlag 1993. Press, 2001.
- [13] **Rose, David, A.**, The College Math Journal, vol. 22, No.2 March 1991.
- [14] **Rudin, W.**, *Principles of mathematical analysis*, McGraw Hill third edition 1976
- [15] **Rudin W.**, *Real and Complex Analysis*, third edition, McGraw-Hill, 1987.
- [16] **Salas S. and Hille E.**, *Calculus One and Several Variables*, Wiley 1990.
- [17] **Tierney John**, *Calculus and Analytic Geometry*, fourth edition, Allyn and Bacon, Boston, 1969.

# Index

- $\cap$ , 21
- $\cup$ , 21
- $\nabla^2$ , 593
- $n^{th}$  term test, 266
  
- absolute convergence, 263
- adjugate, 443, 455
- alternating series test, 267
- amplitude, 64, 142
- annuity
  - ordinary, 31
- antiderivatives, 179
- arc length, 375
- Archimedian property, 28
- area of a cone, 68
- arithmetic mean, 506
- augmented matrix, 38
  
- balance of momentum, 606
- bezier curves, 359
- binomial series, 281
- binomial theorem, 32
  
- capacitance, 394
- carbon dating, 222
- Cartesian coordinates, 294
- catenary, 216
- Cauchy, 88
- Cauchy condensation test, 265
- Cauchy mean value theorem, 133
- Cauchy product, 272
- Cauchy Schwartz inequality, 311
- Cauchy Schwarz, 298
- Cauchy sequence, 106, 341
- Cauchy stress, 607
- center of mass, 553
- central force field, 400
- centrifugal acceleration, 432
- centripetal acceleration, 432
- centripetal force, 400
- chain rule, 485
- change of variables formula, 540
- characteristic equation, 511
  
- circular functions, 139
- circular shells, 211
- closed set, 302
- coefficient of thermal conductivity, 494
- cofactor, 439, 453
- column rank, 444, 456
- comparison test, 263
- completeness, 42
- completeness, 107
- completeness axiom, 42
- completing the square, 44
- complex conjugate, 73
- complex numbers, 73
- component, 307
- concave down, 135
- concave up, 135
- conformable, 414
- conservation of linear momentum, 364
- conservation of mass, 605
- conservative, 622
- constitutive laws, 610
- contented set, 566
- continuous function, 88, 334
- Coordinates, 293
- Coriolis acceleration, 432
- Coriolis acceleration
  - earth, 434
- Coriolis force, 400, 432
- Cramer's rule, 444, 456
- critical points, 509
- curl, 593
- curvature, 381
- cycloid, 627
  
- Darboux, 240
- Darboux integral, 240
- deformation gradient, 606
- dense, 29
- density of rationals, 29
- derivative
  - intermediate value property, 134
  - mean value theorem, 132



- derivative of a function, 120, 352
- determinant, 449
  - product, 452
  - transpose, 450
- difference quotient, 120, 352
- differentiability and continuity, 121
- differentiable, 489
- differentiable matrix, 428
- differential equation, 182
- differentiation rules, 122, 355
- directrix, 70, 321
- Dirichlet function, 81
- Dirichlet test, 267
- discriminant, 44
- distance, 53
- distance formula, 296
- divergence, 593
- divergence theorem, 598
- domain, 81
- donut, 583
- dot product, 311
  
- eigenvalue, 505
- eigenvalues, 511
- Einstein summation convention, 331
- ellipse, 70
- equality of mixed partial derivatives, 475
- Euclidean algorithm, 30
- Euler method, 371
- Eulerian coordinates, 606
- extreme value theorem, 96
  
- Fermat's principle, 173
- Fibonacci sequence, 85
- Fick's law, 494, 611
- field axioms, 16
- first order linear differential equations, 250
- focus, 70
- force field, 378, 400
- Foucault pendulum, 434
- Fourier law of heat conduction, 494
- Fredholm alternative, 426
- frequency, 142
- frustum of a cone, 216
- function
  - even, 131
  - odd, 131
  - uniformly continuous, 113
- fundamental matrix, 465
- fundamental theorem of algebra, 629
- fundamental theorem of arithmetic, 34
- fundamental theorem of calculus, 238
- future value of an annuity, 31
  
- Gauss's theorem, 598
- geometric mean, 506
- geometric series, 261
- gradient vector, 493
- Grammian, 461
- greatest common divisor, 33
- greatest lower bound, 42
- grids, 558
  
- hanging chain, 215
- Heine Borel, 113
- Hessian matrix, 510
- Holder's inequality, 175
- Hooke's law, 387
- hyperbola, 71
  
- implicit differentiation, 150
- implicit function theorem, 150
- improper integrals, 252
- inconsistent, 39
- indefinite integral, 179
- inductance, 394
- inflection point, 137
- initial value problem, 179, 182
- inner product, 311
- integral, 179
- integral test, 275
- integrand, 179
- integration by parts, 190, 248
- interior point, 302
- intermediate value theorem, 94, 112
- intervals, 22
- inverse function theorem, 150
- inverse image, 83
- inverses and determinants, 442, 454
- isocles triangle, 55
  
- Jacobian determinant, 540
- Jordan content, 566
- Jordan set, 566
- joule, 316
  
- Kepler's first law, 402
- Kepler's laws, 400
- Kepler's third law, 405
- kilogram, 327
- kinetic energy, 363
- Kroneker delta, 330

- Lagrange multipliers, 503
- Lagrange remainder, 258
- Lagrangian coordinates, 606
- Laplace expansion, 439, 453
- law of cosines, 54
- law of sines, 64
- least squares regression, 473
- least upper bound, 42
- Lebesgue's theorem, 571
- length of smooth curve, 376
- limit comparison test, 264
- limit of a function, 97, 336, 351
- limit point, 475
- line integral, 379
- linear combination, 423, 444, 451, 456
- linear momentum, 364
- Lipschitz, 114, 335, 343
- Lipschitz condition, 367
- lizards
  - surface area, 579
- local maximum, 126
- local minimum, 126
- logarithmic differentiation, 155
- lower sum, 559
  
- main diagonal, 440
- mass ballance, 605
- material coordinates, 606
- mathematical induction, 27
- matrix, 411
  - inverse, 417
  - left inverse, 455
  - lower triangular, 440, 456
  - right inverse, 455
  - upper triangular, 440, 456
- max. min.theorem, 96
- mean value theorem
  - for integrals, 241
- Merten's theorem, 273
- metric tensor, 461
- minor, 439, 453
- Mobius band, 615
- motion, 606
- moving coordinate system, 429
  - acceleration , 431
- multi-index, 335
- multinomial expansion, 44
  
- natural logarithms, 147
- Navier, 612
- nested interval lemma, 95
  
- Newton, 308
  - second law, 359
- Newton Raphson method, 175
- nonremovable discontinuity, 87
  
- Ohm's law, 394
- open set, 302
- order axioms, 20
- ordinary differential equations
  - existence and uniqueness, 370
  - global existence, 370
  - linear systems, 463
  - local existence, 498
- orientation, 378
- oriented curve, 378
- oscillation
  - critically damped, 394
  - over damped, 394
  - underdamped, 394
- osculating plane, 381
  
- p series, 265
- parabola, 69
- parametrization, 375
- partial derivative, 471
- partial fractions expansion, 199
- partial summation formula, 267
- particular solution, 466
- partition, 228
- Pathagorean theorem, 48
- Peano existence theorem, 499
- permutation symbol, 330
- permutations and combinations, 43
- phase shift, 64, 142
- Picard iteration, 368
- Piola Kirchhoff stress, 610
- polar form complex number, 74
- power series, 275
- precession of a top, 551
- present value of an annuity, 31
- prime number, 33
- principal normal, 381
- product rule, 122
  - cross product, 355
  - dot product, 355
  - matrices, 428
- profit, 127
- properties of integral
  - properties, 237
  
- quadratic formula, 44
- quotient rule, 122

- Raabe's test, 275
- radius of curvature, 381
- range, 81
- rank of a matrix, 444, 456
- ratio test, 267
- rational function, 83
- rational root theorem, 34
- rational function of cosines and sines, 202
- real numbers, 15
- recurrence relation, 85
- recursively defined sequence, 112
- recursively defined sequence, 85
- refinement of grids, 558
- regression line, 425
- relatively prime, 33
- removable discontinuity, 87
- resistance, 394
- resultant, 307
- Revenue, 127
- Riemann criterion, 230
- Riemann criterion, 560
- Riemann integrable, 229
- Riemann integral, 560
- rigid body motion, 555
- Rolle's theorem, 132
- root test, 268
- rot, 593
- row operations, 38
- row rank, 444, 456
- 
- saddle point, 510
- scalar field, 593
- scalar potential, 622
- scalar product, 311
- scalars, 294
- separable differential equations, 220
- sequence of partial sums, 261
- sequences, 84
- sequential compactness, 113
- sequentially compact set, 347
- Simpson's rule, 246, 251
- slope, 51
- smooth curve, 375
- Snell's law, 173
- spacial coordinates, 606
- span, 423, 451
- spherical coordinates, 486
- squeezing theorem, 100
- squeezing theorem, 105
- Stokes, 612
- subtend, 56
- 
- superposition, 464
- 
- Taylor series, 275
- Taylor's formula, 250, 257, 508
- torsion, 382
- torus, 583
- trapezoidal rule, 245, 251
- triangle inequality, 24, 299, 312
- trichotomy, 21
- trigonometric functions, 54
- 
- uniformly continuous, 113
- unit tangent vector, 381
- upper and lower sums, 228
- upper sum, 559
- 
- variation of constants formula, 465
- vector field, 378, 593
- vector potential, 596
- vector space, 412
- vector valued function
  - derivative, 352
  - integral, 352
- vectors, 305, 412
- volume element, 540, 580
- 
- Weierstrass, 88
- well ordered, 27
- work, 379
- Wronskian, 465