

Contents

Preface	x
1 Introduction	1
1.1 What is analysis?	1
1.2 Why do analysis?	3
2 The natural numbers	14
2.1 The Peano axioms	16
2.2 Addition	27
2.3 Multiplication	33
3 Set theory	37
3.1 Fundamentals	37
3.2 Russell's paradox (Optional)	52
3.3 Functions	55
3.4 Images and inverse images	64
3.5 Cartesian products	70
3.6 Cardinality of sets	78
4 Integers and rationals	85
4.1 The integers	85
4.2 The rationals	93
4.3 Absolute value and exponentiation	99
4.4 Gaps in the rational numbers	104
5 The real numbers	108
5.1 Cauchy sequences	110

5.2	Equivalent Cauchy sequences	115
5.3	The construction of the real numbers	118
5.4	Ordering the reals	128
5.5	The least upper bound property	134
5.6	Real exponentiation, part I	140
6	Limits of sequences	146
6.1	The Extended real number system	154
6.2	Suprema and Infima of sequences	158
6.3	Limsup, Liminf, and limit points	161
6.4	Some standard limits	171
6.5	Subsequences	172
6.6	Real exponentiation, part II	176
7	Series	179
7.1	Finite series	179
7.2	Infinite series	189
7.3	Sums of non-negative numbers	195
7.4	Rearrangement of series	200
7.5	The root and ratio tests	204
8	Infinite sets	208
8.1	Countability	208
8.2	Summation on infinite sets	216
8.3	Uncountable sets	224
8.4	The axiom of choice	228
8.5	Ordered sets	232
9	Continuous functions on \mathbf{R}	243
9.1	Subsets of the real line	244
9.2	The algebra of real-valued functions	251
9.3	Limiting values of functions	254
9.4	Continuous functions	262
9.5	Left and right limits	267
9.6	The maximum principle	270
9.7	The intermediate value theorem	275
9.8	Monotonic functions	277

9.9	Uniform continuity	280
9.10	Limits at infinity	287
10	Differentiation of functions	290
10.1	Local maxima, local minima, and derivatives . . .	297
10.2	Monotone functions and derivatives	300
10.3	Inverse functions and derivatives	302
10.4	L'Hôpital's rule	305
11	The Riemann integral	308
11.1	Partitions	309
11.2	Piecewise constant functions	314
11.3	Upper and lower Riemann integrals	318
11.4	Basic properties of the Riemann integral	323
11.5	Riemann integrability of continuous functions . . .	329
11.6	Riemann integrability of monotone functions . . .	332
11.7	A non-Riemann integrable function	335
11.8	The Riemann-Stieltjes integral	336
11.9	The two fundamental theorems of calculus	340
11.10	Consequences of the fundamental theorems	345
12	Appendix: the basics of mathematical logic	351
12.1	Mathematical statements	352
12.2	Implication	360
12.3	The structure of proofs	366
12.4	Variables and quantifiers	369
12.5	Nested quantifiers	374
12.6	Some examples of proofs and quantifiers	377
12.7	Equality	379
13	Appendix: the decimal system	382
13.1	The decimal representation of natural numbers . . .	383
13.2	The decimal representation of real numbers	387
14	Metric spaces	391
14.1	Some point-set topology of metric spaces	402
14.2	Relative topology	408

14.3	Cauchy sequences and complete metric spaces . . .	410
14.4	Compact metric spaces	415
15	Continuous functions on metric spaces	422
15.1	Continuity and product spaces	425
15.2	Continuity and compactness	429
15.3	Continuity and connectedness	432
15.4	Topological spaces (Optional)	436
16	Uniform convergence	443
16.1	Limiting values of functions	444
16.2	Pointwise convergence and uniform convergence . .	447
16.3	Uniform convergence and continuity	452
16.4	The metric of uniform convergence	456
16.5	Series of functions; the Weierstrass M -test	459
16.6	Uniform convergence and integration	461
16.7	Uniform convergence and derivatives	464
16.8	Uniform approximation by polynomials	467
17	Power series	477
17.1	Formal power series	477
17.2	Real analytic functions	481
17.3	Abel's theorem	487
17.4	Multiplication of power series	490
17.5	The exponential and logarithm functions	493
17.6	A digression on complex numbers	498
17.7	Trigonometric functions	506
18	Fourier series	514
18.1	Periodic functions	515
18.2	Inner products on periodic functions	518
18.3	Trigonometric polynomials	522
18.4	Periodic convolutions	525
18.5	The Fourier and Plancherel theorems	530

19 Several variable differential calculus	537
19.1 Linear transformations	537
19.2 Derivatives in several variable calculus	545
19.3 Partial and directional derivatives	548
19.4 The several variable calculus chain rule	556
19.5 Double derivatives and Clairaut's theorem	560
19.6 The contraction mapping theorem	562
19.7 The inverse function theorem in several variable calculus	566
19.8 The implicit function theorem	571
20 Lebesgue measure	577
20.1 The goal: Lebesgue measure	579
20.2 First attempt: Outer measure	581
20.3 Outer measure is not additive	591
20.4 Measurable sets	594
20.5 Measurable functions	601
21 Lebesgue integration	605
21.1 Simple functions	605
21.2 Integration of non-negative measurable functions	611
21.3 Integration of absolutely integrable functions	620
21.4 Comparison with the Riemann integral	625
21.5 Fubini's theorem	627

Preface

This text originated from the lecture notes I gave teaching the honours undergraduate-level real analysis sequence at the University of California, Los Angeles, in 2003. Among the undergraduates here, real analysis was viewed as being one of the most difficult courses to learn, not only because of the abstract concepts being introduced for the first time (e.g., topology, limits, measurability, etc.), but also because of the level of rigour and proof demanded of the course. Because of this perception of difficulty, one often was faced with the difficult choice of either reducing the level of rigour in the course in order to make it easier, or to maintain strict standards and face the prospect of many undergraduates, even many of the bright and enthusiastic ones, struggle with the course material.

Faced with this dilemma, I tried a somewhat unusual approach to the subject. Typically, an introductory sequence in real analysis assumes that the students are already familiar with the real numbers, with mathematical induction, with elementary calculus, and with the basics of set theory, and then quickly launches into the heart of the subject, for instance beginning with the concept of a limit. Normally, students entering this sequence do indeed have a fair bit of exposure to these prerequisite topics, however in most cases the material was not covered in a thorough manner; for instance, very few students were able to actually *define* a real number, or even an integer, properly, even though they could visualize these numbers intuitively and manipulate them algebraically.

This seemed to me to be a missed opportunity. Real analysis is one of the first subjects (together with linear algebra and abstract algebra) that a student encounters, in which one truly has to grapple with the subtleties of a truly rigorous mathematical proof. As such, the course offers an excellent chance to go back to the foundations of mathematics - and in particular, the construction of the real numbers - and do it properly and thoroughly.

Thus the course was structured as follows. In the first week, I described some well-known “paradoxes” in analysis, in which standard laws of the subject (e.g., interchange of limits and sums, or sums and integrals) were applied in a non-rigorous way to give nonsensical results such as $0 = 1$. This motivated the need to go back to the very beginning of the subject, even to the very definition of the natural numbers, and check all the foundations from scratch. For instance, one of the first homework assignments was to check (using only the Peano axioms) that addition was associative for natural numbers (i.e., that $(a + b) + c = a + (b + c)$ for all natural numbers a, b, c : see Exercise 2.2.1). Thus even in the first week, the students had to write rigorous proofs using mathematical induction. After we had derived all the basic properties of the natural numbers, we then moved on to the integers (initially defined as formal differences of natural numbers); once the students had verified all the basic properties of the integers, we moved on to the rationals (initially defined as formal quotients of integers); and then from there we moved on (via formal limits of Cauchy sequences) to the reals. Around the same time, we covered the basics of set theory, for instance demonstrating the uncountability of the reals. Only then (after about ten lectures) did we begin what one normally considers the heart of undergraduate real analysis - limits, continuity, differentiability, and so forth.

The response to this format was quite interesting. In the first few weeks, the students found the material very easy on a conceptual level - as we were dealing only with the basic properties of the standard number systems - but very challenging on an intellectual level, as one was analyzing these number systems from a foundational viewpoint for the first time, in order to rigorously derive

the more advanced facts about these number systems from the more primitive ones. One student told me how difficult it was to explain to his friends in the non-honours real analysis sequence (a) why he was still learning how to show why all rational numbers are either positive, negative, or zero (Exercise 4.2.4), while the non-honours sequence was already distinguishing absolutely convergent and conditionally convergent series, and (b) why, despite this, he thought his homework was significantly harder than that of his friends. Another student commented to me, quite wryly, that while she could obviously *see* why one could always divide one positive integer q into natural number n to give a quotient a and a remainder r less than q (Exercise 2.3.5), she still had, to her frustration, much difficulty writing down a proof of this fact. (I told her that later in the course she would have to prove statements for which it would not be as obvious to see that the statements were true; she did not seem to be particularly consoled by this.) Nevertheless, these students greatly enjoyed the homework, as when they did persevere and obtain a rigorous proof of an intuitive fact, it solidified the link in their minds between the abstract manipulations of formal mathematics and their informal intuition of mathematics (and of the real world), often in a very satisfying way. By the time they were assigned the task of giving the infamous “epsilon and delta” proofs in real analysis, they had already had so much experience with formalizing intuition, and in discerning the subtleties of mathematical logic (such as the distinction between the “for all” quantifier and the “there exists” quantifier), that the transition to these proofs was fairly smooth, and we were able to cover material both thoroughly and rapidly. By the tenth week, we had caught up with the non-honours class, and the students were verifying the change of variables formula for Riemann-Stieltjes integrals, and showing that piecewise continuous functions were Riemann integrable. By the the conclusion of the sequence in the twentieth week, we had covered (both in lecture and in homework) the convergence theory of Taylor and Fourier series, the inverse and implicit function theorem for continuously differentiable functions of several variables, and established

the dominated convergence theorem for the Lebesgue integral.

In order to cover this much material, many of the key foundational results were left to the student to prove as homework; indeed, this was an essential aspect of the course, as it ensured the students truly appreciated the concepts as they were being introduced. This format has been retained in this text; the majority of the exercises consist of proving lemmas, propositions and theorems in the main text. Indeed, I would strongly recommend that one do as many of these exercises as possible - and this includes those exercises proving “obvious” statements - if one wishes to use this text to learn real analysis; this is not a subject whose subtleties are easily appreciated just from passive reading. Most of the chapter sections have a number of exercises, which are listed at the end of the section.

To the expert mathematician, the pace of this book may seem somewhat slow, especially in early chapters, as there is a heavy emphasis on rigour (except for those discussions explicitly marked “Informal”), and justifying many steps that would ordinarily be quickly passed over as being self-evident. The first few chapters develop (in painful detail) many of the “obvious” properties of the standard number systems, for instance that the sum of two positive real numbers is again positive (Exercise 5.4.1), or that given any two distinct real numbers, one can find rational number between them (Exercise 5.4.4). In these foundational chapters, there is also an emphasis on *non-circularity* - not using later, more advanced results to prove earlier, more primitive ones. In particular, the usual laws of algebra are not used until they are derived (and they have to be derived separately for the natural numbers, integers, rationals, and reals). The reason for this is that it allows the students to learn the art of abstract reasoning, deducing true facts from a limited set of assumptions, in the friendly and intuitive setting of number systems; the payoff for this practice comes later, when one has to utilize the same type of reasoning techniques to grapple with more advanced concepts (e.g., the Lebesgue integral).

The text here evolved from my lecture notes on the subject, and thus is very much oriented towards a pedagogical perspective;

much of the key material is contained inside exercises, and in many cases I have chosen to give a lengthy and tedious, but instructive, proof instead of a slick abstract proof. In more advanced textbooks, the student will see shorter and more conceptually coherent treatments of this material, and with more emphasis on intuition than on rigour; however, I feel it is important to know how to to analysis rigourously and “by hand” first, in order to truly appreciate the more modern, intuitive and abstract approach to analysis that one uses at the graduate level and beyond.

Some of the material in this textbook is somewhat peripheral to the main theme and may be omitted for reasons of time constraints. For instance, as set theory is not as fundamental to analysis as are the number systems, the chapters on set theory (Chapters 3, 8) can be covered more quickly and with substantially less rigour, or be given as reading assignments. Similarly for the appendices on logic and the decimal system. The chapter on Fourier series is also not needed elsewhere in the text and can be omitted.

I am deeply indebted to my students, who over the progression of the real analysis sequence found many corrections and suggestions to the notes, which have been incorporated here. I am also very grateful to the anonymous referees who made several corrections and suggested many important improvements to the text.

Terence Tao

Chapter 1

Introduction

1.1 What is analysis?

This text is an honours-level undergraduate introduction to *real analysis*: the analysis of the real numbers, sequences and series of real numbers, and real-valued functions. This is related to, but is distinct from, *complex analysis*, which concerns the analysis of the complex numbers and complex functions, *harmonic analysis*, which concerns the analysis of harmonics (waves) such as sine waves, and how they synthesize other functions via the Fourier transform, *functional analysis*, which focuses much more heavily on functions (and how they form things like vector spaces), and so forth. *Analysis* is the rigorous study of such objects, with a focus on trying to pin down precisely and accurately the qualitative and quantitative behavior of these objects. Real analysis is the theoretical foundation which underlies *calculus*, which is the collection of computational algorithms which one uses to manipulate functions.

In this text we will be studying many objects which will be familiar to you from freshman calculus: numbers, sequences, series, limits, functions, definite integrals, derivatives, and so forth. You already have a great deal of experience knowing how to *compute* with these objects; however here we will be focused more on the underlying theory for these objects. We will be concerned with questions such as the following:

1. What is a real number? Is there a largest real number? After 0, what is the “next” real number (i.e., what is the smallest positive real number)? Can you cut a real number into pieces infinitely many times? Why does a number such as 2 have a square root, while a number such as -2 does not? If there are infinitely many reals and infinitely many rationals, how come there are “more” real numbers than rational numbers?
2. How do you take the limit of a sequence of real numbers? Which sequences have limits and which ones don't? If you can stop a sequence from escaping to infinity, does this mean that it must eventually settle down and converge? Can you add infinitely many real numbers together and still get a finite real number? Can you add infinitely many rational numbers together and end up with a non-rational number? If you rearrange the elements of an infinite sum, is the sum still the same?
3. What is a function? What does it mean for a function to be continuous? differentiable? integrable? bounded? can you add infinitely many functions together? What about taking limits of sequences of functions? Can you differentiate an infinite series of functions? What about integrating? If a function $f(x)$ takes the value of $f(0) = 3$ when $x = 0$ and $f(1) = 5$ when $x = 1$, does it have to take every intermediate value between 3 and 5 when x goes between 0 and 1? Why?

You may already know how to answer some of these questions from your calculus classes, but most likely these sorts of issues were only of secondary importance to those courses; the emphasis was on getting you to perform computations, such as computing the integral of $x \sin(x^2)$ from $x = 0$ to $x = 1$. But now that you are comfortable with these objects and already know how to do all the computations, we will go back to the theory and try to *really* understand what is going on.

1.2 Why do analysis?

It is a fair question to ask, “why bother?”, when it comes to analysis. There is a certain philosophical satisfaction in knowing *why* things work, but a pragmatic person may argue that one only needs to know *how* things work to do real-life problems. The calculus training you receive in introductory classes is certainly adequate for you to begin solving many problems in physics, chemistry, biology, economics, computer science, finance, engineering, or whatever else you end up doing - and you can certainly use things like the chain rule, L’Hôpital’s rule, or integration by parts without knowing why these rules work, or whether there are any exceptions to these rules. However, one can get into trouble if one applies rules without knowing where they came from and what the limits of their applicability are. Let me give some examples in which several of these familiar rules, if applied blindly without knowledge of the underlying analysis, can lead to disaster.

Example 1.2.1 (Division by zero). This is a very familiar one to you: the cancellation law $ac = bc \implies a = b$ does not work when $c = 0$. For instance, the identity $1 \times 0 = 2 \times 0$ is true, but if one blindly cancels the 0 then one obtains $1 = 2$, which is false. In this case it was obvious that one was dividing by zero; but in other cases it can be more hidden.

Example 1.2.2 (Divergent series). You have probably seen geometric series such as the infinite sum

$$S = 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16} + \dots$$

You have probably seen the following trick to sum this series: if we call the above sum S , then if we multiply both sides by 2, we obtain

$$2S = 2 + 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \dots = 2 + S$$

and hence $S = 2$, so the series sums to 2. However, if you apply the same trick to the series

$$S = 1 + 2 + 4 + 8 + 16 + \dots$$

one gets nonsensical results:

$$2S = 2 + 4 + 8 + 16 + \dots = S - 1 \implies S = -1.$$

So the same reasoning that shows that $1 + \frac{1}{2} + \frac{1}{4} + \dots = 2$ also gives that $1 + 2 + 4 + 8 + \dots = -1$. Why is it that we trust the first equation but not the second? A similar example arises with the series

$$S = 1 - 1 + 1 - 1 + 1 - 1 + \dots;$$

we can write

$$S = 1 - (1 - 1 + 1 - 1 + \dots) = 1 - S$$

and hence that $S = 1/2$; or instead we can write

$$S = (1 - 1) + (1 - 1) + (1 - 1) + \dots = 0 + 0 + \dots$$

and hence that $S = 0$; or instead we can write

$$S = 1 + (-1 + 1) + (-1 + 1) + \dots = 1 + 0 + 0 + \dots$$

and hence that $S = 1$. Which one is correct? (See Exercise 7.2.1 for an answer.)

Example 1.2.3 (Divergent sequences). Here is a slight variation of the previous example. Let x be a real number, and let L be the limit

$$L = \lim_{n \rightarrow \infty} x^n.$$

Changing variables $n = m + 1$, we have

$$L = \lim_{m+1 \rightarrow \infty} x^{m+1} = \lim_{m+1 \rightarrow \infty} x \times x^m = x \lim_{m+1 \rightarrow \infty} x^m.$$

But if $m + 1 \rightarrow \infty$, then $m \rightarrow \infty$, thus

$$\lim_{m+1 \rightarrow \infty} x^m = \lim_{m \rightarrow \infty} x^m = \lim_{n \rightarrow \infty} x^n = L,$$

and thus

$$xL = L.$$

At this point we could cancel the L 's and conclude that $x = 1$ for an arbitrary real number x , which is absurd. But since we are already aware of the division by zero problem, we could be a little smarter and conclude instead that either $x = 1$, or $L = 0$. In particular we seem to have shown that

$$\lim_{n \rightarrow \infty} x^n = 0 \text{ for all } x \neq 1.$$

But this conclusion is absurd if we apply it to certain values of x , for instance by specializing to the case $x = 2$ we could conclude that the sequence $1, 2, 4, 8, \dots$ converges to zero, and by specializing to the case $x = -1$ we conclude that the sequence $1, -1, 1, -1, \dots$ also converges to zero. These conclusions appear to be absurd; what is the problem with the above argument? (See Exercise 6.2.4 for an answer.)

Example 1.2.4 (Limiting values of functions). Start with the expression $\lim_{x \rightarrow \infty} \sin(x)$, make the change of variable $x = y + \pi$ and recall that $\sin(y + \pi) = -\sin(y)$ to obtain

$$\lim_{x \rightarrow \infty} \sin(x) = \lim_{y + \pi \rightarrow \infty} \sin(y + \pi) = \lim_{y \rightarrow \infty} (-\sin(y)) = -\lim_{y \rightarrow \infty} \sin(y).$$

Since $\lim_{x \rightarrow \infty} \sin(x) = \lim_{y \rightarrow \infty} \sin(y)$ we thus have

$$\lim_{x \rightarrow \infty} \sin(x) = -\lim_{x \rightarrow \infty} \sin(x)$$

and hence

$$\lim_{x \rightarrow \infty} \sin(x) = 0.$$

If we then make the change of variables $x = \pi/2 - z$ and recall that $\sin(\pi/2 - z) = \cos(z)$ we conclude that

$$\lim_{x \rightarrow \infty} \cos(x) = 0.$$

Squaring both of these limits and adding we see that

$$\lim_{x \rightarrow \infty} (\sin^2(x) + \cos^2(x)) = 0^2 + 0^2 = 0.$$

On the other hand, we have $\sin^2(x) + \cos^2(x) = 1$ for all x . Thus we have shown that $1 = 0$! What is the difficulty here?

Example 1.2.5 (Interchanging sums). Consider the following fact of arithmetic. Consider any matrix of numbers, e.g.

$$\begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix}$$

and compute the sums of all the rows and the sums of all the columns, and then total all the row sums and total all the column sums. In both cases you will get the same number - the total sum of all the entries in the matrix:

$$\begin{array}{cccc} \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix} & 6 & & \\ & 15 & & \\ & 24 & & \\ 12 & 15 & 18 & 45 \end{array}$$

To put it another way, if you want to add all the entries in a $m \times n$ matrix together, it doesn't matter whether you sum the rows first or sum the columns first, you end up with the same answer. (Before the invention of computers, accountants and book-keepers would use this fact to guard against making errors when balancing their books.) In series notation, this fact would be expressed as

$$\sum_{i=1}^m \sum_{j=1}^n a_{ij} = \sum_{j=1}^n \sum_{i=1}^m a_{ij},$$

if a_{ij} denoted the entry in the i^{th} row and j^{th} column of the matrix.

Now one might think that this rule should extend easily to infinite series:

$$\sum_{i=1}^{\infty} \sum_{j=1}^{\infty} a_{ij} = \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} a_{ij}.$$

Indeed, if you use infinite series a lot in your work, you will find yourself having to switch summations like this fairly often. Another way of saying this fact is that in an infinite matrix, the sum of the row-totals should equal the sum of the column-totals.

However, despite the reasonableness of this statement, it is actually false! Here is a counterexample:

$$\begin{pmatrix} 1 & 0 & 0 & 0 & \dots \\ -1 & 1 & 0 & 0 & \dots \\ 0 & -1 & 1 & 0 & \dots \\ 0 & 0 & -1 & 1 & \dots \\ 0 & 0 & 0 & -1 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{pmatrix}.$$

If you sum up all the rows, and then add up all the row totals, you get 1; but if you sum up all the columns, and add up all the column totals, you get 0! So, does this mean that summations for infinite series should not be swapped, and that any argument using such a swapping should be distrusted? (See Theorem 8.2.2 for an answer.)

Example 1.2.6 (Interchanging integrals). The interchanging of integrals is a trick which occurs just as commonly as integrating by sums in mathematics. Suppose one wants to compute the volume under a surface $z = f(x, y)$ (let us ignore the limits of integration for the moment). One can do it by slicing parallel to the x -axis: for each fixed value of y , we can compute an area $\int f(x, y) dx$, and then we integrate the area in the y variable to obtain the volume

$$V = \int \int f(x, y) dx dy.$$

Or we could slice parallel to the y -axis to obtain an area $\int f(x, y) dy$, and then integrate in the x -axis to obtain

$$V = \int \int f(x, y) dy dx.$$

This seems to suggest that one should always be able to swap integral signs:

$$\int \int f(x, y) dx dy = \int \int f(x, y) dy dx.$$

And indeed, people swap integral signs all the time, because sometimes one variable is easier to integrate in first than the other. However, just as infinite sums sometimes cannot be swapped, integrals are also sometimes dangerous to swap. An example is with the integrand $e^{-xy} - xye^{-xy}$. Suppose we believe that we can swap the integrals:

$$\int_0^\infty \int_0^1 (e^{-xy} - xye^{-xy}) dy dx = \int_0^1 \int_0^\infty (e^{-xy} - xye^{-xy}) dx dy.$$

Since

$$\int_0^1 (e^{-xy} - xye^{-xy}) dy = ye^{-xy} \Big|_{y=0}^{y=1} = e^{-x},$$

the left-hand side is $\int_0^\infty e^{-x} dx = -e^{-x} \Big|_0^\infty = 1$. But since

$$\int_0^\infty (e^{-xy} - xye^{-xy}) dx = xe^{-xy} \Big|_{x=0}^{x=\infty} = 0,$$

the right-hand side is $\int_0^1 0 dx = 0$. Clearly $1 \neq 0$, so there is an error somewhere; but you won't find one anywhere except in the step where we interchanged the integrals. So how do we know when to trust the interchange of integrals? (See Theorem 21.5.1 for a partial answer.)

Example 1.2.7 (Interchanging limits). Suppose we start with the plausible looking statement

$$\lim_{x \rightarrow 0} \lim_{y \rightarrow 0} \frac{x^2}{x^2 + y^2} = \lim_{y \rightarrow 0} \lim_{x \rightarrow 0} \frac{x^2}{x^2 + y^2}. \quad (1.1)$$

But we have

$$\lim_{y \rightarrow 0} \frac{x^2}{x^2 + y^2} = \frac{x^2}{x^2 + 0^2} = 1,$$

so the left-hand side of (1.1) is 1; on the other hand, we have

$$\lim_{x \rightarrow 0} \frac{x^2}{x^2 + y^2} = \frac{0^2}{0^2 + y^2} = 0,$$

so the right-hand side of (1.1) is 0. Since 1 is clearly not equal to zero, this suggests that interchange of limits is untrustworthy. But are there any other circumstances in which the interchange of limits is legitimate? (See Exercise 15.1.9 for a partial answer.)

Example 1.2.8 (Interchanging limits, again). Start with the plausible looking statement

$$\lim_{x \rightarrow 1^-} \lim_{n \rightarrow \infty} x^n = \lim_{n \rightarrow \infty} \lim_{x \rightarrow 1^-} x^n$$

where the notation $x \rightarrow 1^-$ means that x is approaching 1 from the left. When x is to the left of 1, then $\lim_{n \rightarrow \infty} x^n = 0$, and hence the left-hand side is zero. But we also have $\lim_{x \rightarrow 1^-} x^n = 1$ for all n , and so the right-hand side limit is 1. Does this demonstrate that this type of limit interchange is always untrustworthy? (See Proposition 16.3.3 for an answer.)

Example 1.2.9 (Interchanging limits and integrals). For any real number y , we have

$$\int_{-\infty}^{\infty} \frac{1}{1 + (x - y)^2} dx = \arctan(x - y)|_{x=-\infty}^{\infty} = \frac{\pi}{2} - \left(-\frac{\pi}{2}\right) = \pi.$$

Taking limits as $y \rightarrow \infty$, we should obtain

$$\int_{-\infty}^{\infty} \lim_{y \rightarrow \infty} \frac{1}{1 + (x - y)^2} dx = \lim_{y \rightarrow \infty} \int_{-\infty}^{\infty} \frac{1}{1 + (x - y)^2} dx = \pi.$$

But for every x , we have $\lim_{y \rightarrow \infty} \frac{1}{1 + (x - y)^2} = 0$. So we seem to have concluded that $0 = \pi$. What was the problem with the above argument? Should one abandon the (very useful) technique of interchanging limits and integrals? (See Theorem 16.6.1 for a partial answer.)

Example 1.2.10 (Interchanging limits and derivatives). Observe that if $\varepsilon > 0$, then

$$\frac{d}{dx} \left(\frac{x^3}{\varepsilon^2 + x^2} \right) = \frac{3x^2(\varepsilon^2 + x^2) - 2x^4}{(\varepsilon^2 + x^2)^2}$$

and in particular that

$$\frac{d}{dx} \left(\frac{x^3}{\varepsilon^2 + x^2} \right) \Big|_{x=0} = 0.$$

Taking limits as $\varepsilon \rightarrow 0$, one might then expect that

$$\frac{d}{dx} \left(\frac{x^3}{0+x^2} \right) \Big|_{x=0} = 0.$$

But the right-hand side is $\frac{d}{dx}x = 1$. Does this mean that it is always illegitimate to interchange limits and derivatives? (see Theorem 16.7.1 for an answer.)

Example 1.2.11 (Interchanging derivatives). Let¹ $f(x, y) := \frac{xy^3}{x^2+y^2}$. A common manoeuvre in analysis is to interchange two partial derivatives, thus one expects

$$\frac{\partial^2 f}{\partial x \partial y}(0, 0) = \frac{\partial^2 f}{\partial y \partial x}(0, 0).$$

But from the quotient rule we have

$$\frac{\partial f}{\partial y}(x, y) = \frac{3xy^2}{x^2+y^2} - \frac{2xy^4}{(x^2+y^2)^2}$$

and in particular

$$\frac{\partial f}{\partial y}(x, 0) = \frac{0}{x^2} - \frac{0}{x^4} = 0.$$

Thus

$$\frac{\partial^2 f}{\partial x \partial y}(0, 0) = 0.$$

On the other hand, from the quotient rule again we have

$$\frac{\partial f}{\partial x}(x, y) = \frac{y^3}{x^2+y^2} - \frac{2x^2y^3}{(x^2+y^2)^2}$$

and hence

$$\frac{\partial f}{\partial x}(0, y) = \frac{y^3}{y^2} - \frac{0}{y^4} = y.$$

¹One might object that this function is not defined at $(x, y) = (0, 0)$, but if we set $f(0, 0) := (0, 0)$ then this function becomes continuous and differentiable for all (x, y) , and in fact both partial derivatives $\frac{\partial f}{\partial x}$, $\frac{\partial f}{\partial y}$ are also continuous and differentiable for all (x, y) !

Thus

$$\frac{\partial^2 f}{\partial x \partial y}(0, 0) = 1.$$

Since $1 \neq 0$, we thus seem to have shown that interchange of derivatives is untrustworthy. But are there any other circumstances in which the interchange of derivatives is legitimate? (See Theorem 19.5.4 and Exercise 19.5.1 for some answers.)

Example 1.2.12 (L'Hôpital's rule). We are all familiar with the beautifully simple L'Hôpital's rule

$$\lim_{x \rightarrow x_0} \frac{f(x)}{g(x)} = \lim_{x \rightarrow x_0} \frac{f'(x)}{g'(x)},$$

but one can still get led to incorrect conclusions if one applies it incorrectly. For instance, applying it to $f(x) := x$, $g(x) := 1 + x$, and $x_0 := 0$ we would obtain

$$\lim_{x \rightarrow 0} \frac{x}{1+x} = \lim_{x \rightarrow 0} \frac{1}{1} = 1,$$

but this is the incorrect answer, since $\lim_{x \rightarrow 0} \frac{x}{1+x} = \frac{0}{1+0} = 0$. Of course, all that is going on here is that L'Hôpital's rule is only applicable when both $f(x)$ and $g(x)$ go to zero as $x \rightarrow x_0$, a condition which was violated in the above example. But even when $f(x)$ and $g(x)$ do go to zero as $x \rightarrow x_0$ there is still a possibility for an incorrect conclusion. For instance, consider the limit

$$\lim_{x \rightarrow 0} \frac{x^2 \sin(1/x^4)}{x}.$$

Both numerator and denominator go to zero as $x \rightarrow 0$, so it seems pretty safe to apply L'Hôpital's rule, to obtain

$$\begin{aligned} \lim_{x \rightarrow 0} \frac{x^2 \sin(1/x^4)}{x} &= \lim_{x \rightarrow 0} \frac{2x \sin(1/x^4) - 4x^{-2} \cos(1/x^4)}{1} \\ &= \lim_{x \rightarrow 0} 2x \sin(1/x^4) - \lim_{x \rightarrow 0} 4x^{-2} \cos(1/x^4). \end{aligned}$$

The first limit converges to zero by the squeeze test (since $2x \sin(1/x^4)$ is bounded above by $2|x|$ and below by $-2|x|$, both of which go to

zero at 0). But the second limit is divergent (because x^{-2} goes to infinity as $x \rightarrow 0$, and $\cos(1/x^4)$ does not go to zero). So the limit $\lim_{x \rightarrow 0} \frac{2x \sin(1/x^4) - 4x^{-2} \cos(1/x^4)}{1}$ diverges. One might then conclude using L'Hôpital's rule that $\lim_{x \rightarrow 0} \frac{x^2 \sin(1/x^4)}{x}$ also diverges; however we can clearly rewrite this limit as $\lim_{x \rightarrow 0} x \sin(1/x^4)$, which goes to zero when $x \rightarrow 0$ by the squeeze test again. This does not show that L'Hôpital's rule is untrustworthy (indeed, it is quite rigorous; see Section 10.4), but it still requires some care when applied.

Example 1.2.13 (Limits and lengths). When you learn about integration and how it relates to the area under a curve, you were probably presented with some picture in which the area under the curve was approximated by a bunch of rectangles, whose area was given by a Riemann sum, and then one somehow “took limits” to replace that Riemann sum with an integral, which then presumably matched the actual area under the curve. Perhaps a little later, you learnt how to compute the length of a curve by a similar method - approximate the curve by a bunch of line segments, compute the length of all the line segments, then take limits again to see what you get.

However, it should come as no surprise by now that this approach also can lead to nonsense if used incorrectly. Consider the right-angled triangle with vertices $(0, 0)$, $(1, 0)$, and $(0, 1)$, and suppose we wanted to compute the length of the hypotenuse of this triangle. Pythagoras' theorem tells us that this hypotenuse has length $\sqrt{2}$, but suppose for some reason that we did not know about Pythagoras' theorem, and wanted to compute the length using calculus methods. Well, one way to do so is to approximate the hypotenuse by horizontal and vertical edges. Pick a large number N , and approximate the hypotenuse by a “staircase” consisting of N horizontal edges of equal length, alternating with N vertical edges of equal length. Clearly these edges all have length $1/N$, so the total length of the staircase is $2N/N = 2$. If one takes limits as N goes to infinity, the staircase clearly approaches the hypotenuse, and so in the limit we should get the length of the

hypotenuse. However, as $N \rightarrow \infty$, the limit of $2N/N$ is 2, not $\sqrt{2}$, so we have an incorrect value for the length of the hypotenuse. How did this happen?

The analysis you learn in this text will help you resolve these questions, and will let you know when these rules (and others) are justified, and when they are illegal, thus separating the useful applications of these rules from the nonsense. Thus they can prevent you from making mistakes, and can help you place these rules in a wider context. Moreover, as you learn analysis you will develop an “analytical way of thinking”, which will help you whenever you come into contact with any new rules of mathematics, or when dealing with situations which are not quite covered by the standard rules (e.g., what if your functions are complex-valued instead of real-valued? What if you are working on the sphere instead of the plane? What if your functions are not continuous, but are instead things like square waves and delta functions? What if your functions, or limits of integration, or limits of summation, are occasionally infinite?). You will develop a sense of *why* a rule in mathematics (e.g., the chain rule) works, how to adapt it to new situations, and what its limitations (if any) are; this will allow you to apply the mathematics you have already learnt more confidently and correctly.

Chapter 2

Starting at the beginning: the natural numbers

In this text, we will want to go over much of the material you have learnt in high school and in elementary calculus classes, but to do so as rigorously as possible. To do so we will have to begin at the very basics - indeed, we will go back to the concept of *numbers* and what their properties are. Of course, you have dealt with numbers for over ten years and you know very well how to manipulate the rules of algebra to simplify any expression involving numbers, but we will now turn to a more fundamental issue, which is *why* the rules of algebra work... for instance, why is it true that $a(b + c)$ is equal to $ab + ac$ for any three numbers a, b, c ? This is not an arbitrary choice of rule; it can be proven from more primitive, and more fundamental, properties of the number system. This will teach you a new skill - how to prove complicated properties from simpler ones. You will find that even though a statement may be “obvious”, it may not be easy to prove; the material here will give you plenty of practice in doing so, and in the process will lead you to think about *why* an obvious statement really is obvious. One skill in particular that you will pick up here is the use of *mathematical induction*, which is a basic tool in proving things in many areas of mathematics.

So in the first few chapters we will re-acquaint you with various number systems that are used in real analysis. In increasing order of sophistication, they are the *natural numbers* \mathbf{N} ; the *integers*

\mathbf{Z} ; the *rational numbers* \mathbf{Q} , and the *real numbers* \mathbf{R} . (There are other number systems such as the *complex numbers* \mathbf{C} , but we will not study them until Section 17.6.) The natural numbers $\{0, 1, 2, \dots\}$ are the most primitive of the number systems, but they are used to build the integers, which in turn are used to build the rationals, which in turn build the real numbers. Thus to begin at the very beginning, we must look at the natural numbers. We will consider the following question: how does one actually *define* the natural numbers? (This is a very different question as to how to *use* the natural numbers, which is something you of course know how to do very well. It's like the difference between knowing how to use, say, a computer, versus knowing how to *build* that computer.)

This question is more difficult to answer than it looks. The basic problem is that you have used the natural numbers for so long that they are embedded deeply into your mathematical thinking, and you can make various implicit assumptions about these numbers (e.g., that $a + b$ is always equal to $b + a$) without even thinking; it is difficult to let go for a moment and try to inspect this number system as if it is the first time you have seen it. So in what follows I will have to ask you to perform a rather difficult task: try to set aside, for the moment, everything you know about the natural numbers; forget that you know how to count, to add, to multiply, to manipulate the rules of algebra, etc. We will try to introduce these concepts one at a time and try to identify explicitly what our assumptions are as we go along - and not allow ourselves to use more "advanced" tricks - such as the rules of algebra - until we have actually proven them. This may seem like an irritating constraint, especially as we will spend a lot of time proving statements which are "obvious", but it is necessary to do this suspension of known facts to avoid *circularity* (e.g., using an advanced fact to prove a more elementary fact, and then later using the elementary fact to prove the advanced fact). Also, it is an excellent exercise for really affirming the foundations of your mathematical knowledge, and practicing your proofs and abstract thinking here will be invaluable when we move on to more advanced concepts, such as real numbers, then functions, then se-

quences and series, then differentials and integrals, and so forth. In short, the results here may seem trivial, but the journey is much more important than the destination, for now. (Once the number systems are constructed properly, we can resume using the laws of algebra etc. without having to rederive them each time.)

We will also forget that we know the decimal system, which of course is an extremely convenient way to manipulate numbers, but it is not something which is fundamental to what numbers are. (For instance, one could use an octal or binary system instead of the decimal system, or even the Roman numeral system, and still get exactly the same set of numbers.) Besides, if one tries to fully explain what the decimal number system is, it isn't as natural as you might think. Why is 00423 the same number as 423, but 32400 isn't the same number as 324? How come $123.4444\dots$ is a real number, but $\dots 444.321$ isn't? And why do we have to do all this carrying of digits when adding or multiplying? Why is $0.999\dots$ the same number as 1? What is the smallest positive real number? Isn't it just $0.00\dots 001$? So to set aside these problems, we will not try to assume any knowledge of the decimal system (though we will of course still refer to numbers by their familiar names such as 1,2,3, etc. instead of using other notation such as I,II,III or 0++, (0++)++, ((0++)++)++ (see below) so as not to be needlessly artificial). For completeness, we review the decimal system in an Appendix (§13).

2.1 The Peano axioms

We now present one standard way to define the natural numbers, in terms of the *Peano axioms*, which were first laid out by Giuseppe Peano (1858–1932). This is not the only way to define the natural numbers. For instance, another approach is to talk about the cardinality of finite sets, for instance one could take a set of five elements and define 5 to be the number of elements in that set. We shall discuss this alternate approach in Section 3.6. However, we shall stick with the Peano axiomatic approach for now.

How are we to define what the natural numbers are? Informally, we could say

Definition 2.1.1. (Informal) A *natural number* is any element of the set

$$\mathbf{N} := \{0, 1, 2, 3, 4, \dots\},$$

which is the set of all the numbers created by starting with at 0 and then counting forward indefinitely. We call \mathbf{N} the *set of natural numbers*.

Remark 2.1.2. In some texts the natural numbers start at 1 instead of 0, but this is a matter of notational convention more than anything else. In this text we shall refer to the set $\{1, 2, 3, \dots\}$ as the *positive integers* \mathbf{Z}^+ rather than the natural numbers. Natural numbers are sometimes also known as *whole numbers*.

In a sense, this definition solves the problem of what the natural numbers are: a natural number is any element of the set¹ \mathbf{N} . However, it is not really that satisfactory, because it begs the question of what \mathbf{N} is. This definition of “start at 0 and count indefinitely” seems like an intuitive enough definition of \mathbf{N} , but it is not entirely acceptable, because it leaves many questions unanswered. For instance: how do we know we can keep counting indefinitely, without cycling back to 0? Also, how do you perform operations such as addition, multiplication, or exponentiation?

We can answer the latter question first: we can define complicated operations in terms of simpler operations. Exponentiation is nothing more than repeated multiplication: 5^3 is nothing more than three fives multiplied together. Multiplication is nothing more than repeated addition; 5×3 is nothing more than three fives added together. (Subtraction and division will not be covered here, because they are not operations which are well-suited to the natural numbers; they will have to wait for the integers

¹Strictly speaking, there is another problem with this informal definition: we have not yet defined what a “set” is, or what “element of” is. Thus for the rest of this chapter we shall avoid mention of sets and their elements as much as possible, except in informal discussion.

and rationals, respectively.) And addition? It is nothing more than the repeated operation of *counting forward*, or *incrementing*. If you add three to five, what you are doing is incrementing five three times. On the other hand, incrementing seems to be a pretty fundamental operation, not reducible to any simpler operation; indeed, it is the first operation one learns on numbers, even before learning to add.

Thus, to define the natural numbers, we will use two fundamental concepts: the zero number 0, and the increment operation. In deference to modern computer languages, we will use $n++$ to denote the increment or *successor* of n , thus for instance $3++ = 4$, $(3++)++ = 5$, etc. This is slightly different usage from that in computer languages such as C , where $n++$ actually *redefines* the value of n to be its successor; however in mathematics we try not to define a variable more than once in any given setting, as it can often lead to confusion; many of the statements which were true for the old value of the variable can now become false, and vice versa.

So, it seems like we want to say that \mathbf{N} consists of 0 and everything which can be obtained from 0 by incrementing: \mathbf{N} should consist of the objects

$$0, 0++, (0++)++, ((0++)++)++, \text{ etc.}$$

If we start writing down what this means about the natural numbers, we thus see that we should have the following axioms concerning 0 and the increment operation $++$:

Axiom 2.1. *0 is a natural number.*

Axiom 2.2. *If n is a natural number, then $n++$ is also a natural number.*

Thus for instance, from Axiom 2.1 and two applications of Axiom 2.2, we see that $(0++)++$ is a natural number. Of course, this notation will begin to get unwieldy, so we adopt a convention to write these numbers in more familiar notation:

Definition 2.1.3. We define 1 to be the number $0++$, 2 to be the number $(0++)++$, 3 to be the number $((0++)++)++$, etc. (In other words, $1 := 0++$, $2 := 1++$, $3 := 2++$, etc. In this text I use “ $x := y$ ” to denote the statement that x is *defined* to equal y .)

Thus for instance, we have

Proposition 2.1.4. *3 is a natural number.*

Proof. By Axiom 2.1, 0 is a natural number. By Axiom 2.2, $0++ = 1$ is a natural number. By Axiom 2.2 again, $1++ = 2$ is a natural number. By Axiom 2.2 again, $2++ = 3$ is a natural number. \square

It may seem that this is enough to describe the natural numbers. However, we have not pinned down completely the behavior of \mathbf{N} :

Example 2.1.5. Consider a number system which consists of the numbers 0, 1, 2, 3, in which the increment operation wraps back from 3 to 0. More precisely $0++$ is equal to 1, $1++$ is equal to 2, $2++$ is equal to 3, but $3++$ is equal to 0 (and also equal to 4, by definition of 4). This type of thing actually happens in real life, when one uses a computer to try to store a natural number: if one starts at 0 and performs the increment operation repeatedly, eventually the computer will overflow its memory and the number will wrap around back to 0 (though this may take quite a large number of incrementation operations, such as 65,536). Note that this type of number system obeys Axiom 2.1 and Axiom 2.2, even though it clearly does not correspond to what we intuitively believe the natural numbers to be like.

To prevent this sort of “wrap-around issue” we will impose another axiom:

Axiom 2.3. *0 is not the successor of any natural number; i.e., we have $n++ \neq 0$ for every natural number n .*

Now we can show that certain types of wrap-around do not occur: for instance we can now rule out the type of behavior in Example 2.1.5 using

Proposition 2.1.6. *4 is not equal to 0.*

Don't laugh! Because of the way we have defined 4 - it is the increment of the increment of the increment of the increment of 0 - it is not necessarily true *a priori* that this number is not the same as zero, even if it is "obvious". ("a priori" is Latin for "beforehand" - it refers to what one already knows or assumes to be true before one begins a proof or argument. The opposite is "a posteriori" - what one knows to be true after the proof or argument is concluded). Note for instance that in Example 2.1.5, 4 was indeed equal to 0, and that in a standard two-byte computer representation of a natural number, for instance, 65536 is equal to 0 (using our definition of 65536 as equal to 0 incremented sixty-five thousand, five hundred and thirty-six times).

Proof. By definition, $4 = 3++$. By Axioms 2.1 and 2.2, 3 is a natural number. Thus by Axiom 2.3, $3++ \neq 0$, i.e., $4 \neq 0$. \square

However, even with our new axiom, it is still possible that our number system behaves in other pathological ways:

Example 2.1.7. Consider a number system consisting of five numbers 0,1,2,3,4, in which the increment operation hits a "ceiling" at 4. More precisely, suppose that $0++ = 1$, $1++ = 2$, $2++ = 3$, $3++ = 4$, but $4++ = 4$ (or in other words that $5 = 4$, and hence $6 = 4$, $7 = 4$, etc.). This does not contradict Axioms 2.1,2.2,2.3. Another number system with a similar problem is one in which incrementation wraps around, but not to zero, e.g. suppose that $4++ = 1$ (so that $5 = 1$, then $6 = 2$, etc.).

There are many ways to prohibit the above types of behavior from happening, but one of the simplest is to assume the following axiom:

Axiom 2.4. *Different natural numbers must have different successors; i.e., if n, m are natural numbers and $n \neq m$, then $n++ \neq m++$. Equivalently², if $n++ = m++$, then we must have $n = m$.*

²This is an example of reformulating an implication using its *contrapositive*; see Section 12.2 for more details.

Thus, for instance, we have

Proposition 2.1.8. *6 is not equal to 2.*

Proof. Suppose for contradiction that $6 = 2$. Then $5++ = 1++$, so by Axiom 2.4 we have $5 = 1$, so that $4++ = 0++$. By Axiom 2.4 again we then have $4 = 0$, which contradicts our previous proposition. \square

As one can see from this proposition, it now looks like we can keep all of the natural numbers distinct from each other. There is however still one more problem: while the axioms (particularly Axioms 2.1 and 2.2) allow us to confirm that $0, 1, 2, 3, \dots$ are distinct elements of \mathbf{N} , there is the problem that there may be other “rogue” elements in our number system which are not of this form:

Example 2.1.9. (Informal) Suppose that our number system \mathbf{N} consisted of the following collection of integers and half-integers:

$$\mathbf{N} := \{0, 0.5, 1, 1.5, 2, 2.5, 3, 3.5, \dots\}.$$

(This example is marked “informal” since we are using real numbers, which we’re not supposed to use yet.) One can check that Axioms 2.1-2.4 are still satisfied for this set.

What we want is some axiom which says that the only numbers in \mathbf{N} are those which can be obtained from 0 and the increment operation - in order to exclude elements such as 0.5. But it is difficult to quantify what we mean by “can be obtained from” without already using the natural numbers, which we are trying to define. Fortunately, there is an ingenious solution to try to capture this fact:

Axiom 2.5 (Principle of mathematical induction). *Let $P(n)$ be any property pertaining to a natural number n . Suppose that $P(0)$ is true, and suppose that whenever $P(n)$ is true, $P(n++)$ is also true. Then $P(n)$ is true for every natural number n .*

Remark 2.1.10. We are a little vague on what “property” means at this point, but some possible examples of $P(n)$ might be “ n is even”; “ n is equal to 3”; “ n solves the equation $(n + 1)^2 = n^2 + 2n + 1$ ”; and so forth. Of course we haven’t defined many of these concepts yet, but when we do, Axiom 2.5 will apply to these properties. (A logical remark: Because this axiom refers not just to *variables*, but also *properties*, it is of a different nature than the other four axioms; indeed, Axiom 2.5 should technically be called an *axiom schema* rather than an *axiom* - it is a template for producing an (infinite) number of axioms, rather than being a single axiom in its own right. To discuss this distinction further is far beyond the scope of this text, though, and falls in the realm of logic.)

The informal intuition behind this axiom is the following. Suppose $P(n)$ is such that $P(0)$ is true, and such that whenever $P(n)$ is true, then $P(n++)$ is true. Then since $P(0)$ is true, $P(0++) = P(1)$ is true. Since $P(1)$ is true, $P(1++) = P(2)$ is true. Repeating this indefinitely, we see that $P(0)$, $P(1)$, $P(2)$, $P(3)$, etc. are all true - however this line of reasoning will never let us conclude that $P(0.5)$, for instance, is true. Thus Axiom 2.5 should not hold for number systems which contain “unnecessary” elements such as 0.5. (Indeed, one can give a “proof” of this fact. Apply Axiom 2.5 to the property $P(n) = n$ “is not a half-integer”, i.e., an integer plus 0.5. Then $P(0)$ is true, and if $P(n)$ is true, then $P(n++)$ is true. Thus Axiom 2.5 asserts that $P(n)$ is true for all natural numbers n , i.e., no natural number can be a half-integer. In particular, 0.5 cannot be a natural number. This “proof” is not quite genuine, because we have not defined such notions as “integer”, “half-integer”, and “0.5” yet, but it should give you some idea as to how the principle of induction is supposed to prohibit any numbers other than the “true” natural numbers from appearing in \mathbf{N} .)

The principle of induction gives us a way to prove that a property $P(n)$ is true for every natural number n . Thus in the rest of this text we will see many proofs which have a form like this:

Proposition 2.1.11. *A certain property $P(n)$ is true for every natural number n .*

Proof. We use induction. We first verify the base case $n = 0$, i.e., we prove $P(0)$. (Insert proof of $P(0)$ here). Now suppose inductively that n is a natural number, and $P(n)$ has already been proven. We now prove $P(n++)$. (Insert proof of $P(n++)$, assuming that $P(n)$ is true, here). This closes the induction, and thus $P(n)$ is true for all numbers n . \square

Of course we will not necessarily use the exact template, wording, or order in the above type of proof, but the proofs using induction will generally be something like the above form. There are also some other variants of induction which we shall encounter later, such as backwards induction (Exercise 2.2.6), strong induction (Proposition 2.2.14), and transfinite induction (Lemma 8.5.15).

Axioms 2.1-2.5 are known as the *Peano axioms* for the natural numbers. They are all very plausible, and so we shall make

Assumption 2.6. *(Informal) There exists a number system \mathbf{N} , whose elements we will call natural numbers, for which Axioms 2.1-2.5 are true.*

We will make this assumption a bit more precise once we have laid down our notation for sets and functions in the next chapter.

Remark 2.1.12. We will refer to this number system \mathbf{N} as *the* natural number system. One could of course consider the possibility that there is more than one natural number system, e.g., we could have the Hindu-Arabic number system $\{0, 1, 2, 3, \dots\}$ and the Roman number system $\{O, I, II, III, IV, V, VI, \dots\}$, and if we really wanted to be annoying we could view these number systems as different. But these number systems are clearly equivalent (the technical term is “isomorphic”), because one can create a one-to-one correspondence $0 \leftrightarrow O$, $1 \leftrightarrow I$, $2 \leftrightarrow II$, etc. which maps the zero of the Hindu-Arabic system with the zero of the Roman system, and which is preserved by the increment operation (e.g.,

if 2 corresponds to II , then $2++$ will correspond to $II++$). For a more precise statement of this type of equivalence, see Exercise 3.5.13. Since all versions of the natural number system are equivalent, there is no point in having distinct natural number systems, and we will just use a single natural number system to do mathematics.

We will not prove Assumption 2.6 (though we will eventually include it in our axioms for set theory, see Axiom 3.7), and it will be the only assumption we will ever make about our numbers. A remarkable accomplishment of modern analysis is that just by starting from these five very primitive axioms, and some additional axioms from set theory, we can build all the other number systems, create functions, and do all the algebra and calculus that we are used to.

Remark 2.1.13. (Informal) One interesting feature about the natural numbers is that while each individual natural number is finite, the *set* of natural numbers is infinite; i.e., \mathbf{N} is infinite but consists of individually finite elements. (The whole is greater than any of its parts.) There are no infinite natural numbers; one can even prove this using Axiom 2.5, provided one is comfortable with the notions of finite and infinite. (Clearly 0 is finite. Also, if n is finite, then clearly $n++$ is also finite. Hence by Axiom 2.5, all natural numbers are finite.) So the natural numbers can *approach* infinity, but never actually reach it; infinity is not one of the natural numbers. (There are other number systems which admit “infinite” numbers, such as the cardinals, ordinals, and p-adics, but they do not obey the principle of induction, and in any event are beyond the scope of this text.)

Remark 2.1.14. Note that our definition of the natural numbers is *axiomatic* rather than *constructive*. We have not told you what the natural numbers *are* (so we do not address such questions as what the numbers are made of, are they physical objects, what do they measure, etc.) - we have only listed some things you can do with them (in fact, the only operation we have defined

on them right now is the increment one) and some of the properties that they have. This is how mathematics works - it treats its objects *abstractly*, caring only about what properties the objects have, not what the objects are or what they mean. If one wants to do mathematics, it does not matter whether a natural number means a certain arrangement of beads on an abacus, or a certain organization of bits in a computer's memory, or some more abstract concept with no physical substance; as long as you can increment them, see if two of them are equal, and later on do other arithmetic operations such as add and multiply, they qualify as numbers for mathematical purposes (provided they obey the requisite axioms, of course). It is possible to construct the natural numbers from other mathematical objects - from sets, for instance - but there are multiple ways to construct a working model of the natural numbers, and it is pointless, at least from a mathematician's standpoint, as to argue about which model is the "true" one - as long as it obeys all the axioms and does all the right things, that's good enough to do maths.

Remark 2.1.15. Historically, the realization that numbers could be treated axiomatically is very recent, not much more than a hundred years old. Before then, numbers were generally understood to be inextricably connected to some external concept, such as counting the cardinality of a set, measuring the length of a line segment, or the mass of a physical object, etc. This worked reasonably well, until one was forced to move from one number system to another; for instance, understanding numbers in terms of counting beads, for instance, is great for conceptualizing the numbers 3 and 5, but doesn't work so well for -3 or $1/3$ or $\sqrt{2}$ or $3+4i$; thus each great advance in the theory of numbers - negative numbers, irrational numbers, complex numbers, even the number zero - led to a great deal of unnecessary philosophical anguish. The great discovery of the late nineteenth century was that numbers can be understood abstractly via axioms, without necessarily needing a conceptual model; of course a mathematician can use any of these models when it is convenient, to aid his or her intuition and understanding, but they can also be just as easily discarded when

they begin to get in the way.

One consequence of the axioms is that we can now define sequences *recursively*. Suppose we want to build a sequence a_0, a_1, a_2, \dots by first defining a_0 to be some base value, e.g., $a_0 := c$, and then letting a_1 be some function of a_0 , $a_1 := f_0(a_0)$, a_2 be some function of a_1 , $a_2 := f_1(a_1)$, and so forth - in general, we set $a_{n++} := f_n(a_n)$ for some function f_n from \mathbf{N} to \mathbf{N} . By using all the axioms together we will now conclude that this procedure will give a single value to the sequence element a_n for each natural number n . More precisely³:

Proposition 2.1.16 (Recursive definitions). *Suppose for each natural number n , we have some function $f_n : \mathbf{N} \rightarrow \mathbf{N}$ from the natural numbers to the natural numbers. Let c be a natural number. Then we can assign a unique natural number a_n to each natural number n , such that $a_0 = c$ and $a_{n++} = f_n(a_n)$ for each natural number n .*

Proof. (Informal) We use induction. We first observe that this procedure gives a single value to a_0 , namely c . (None of the other definitions $a_{n++} := f_n(a_n)$ will redefine the value of a_0 , because of Axiom 2.3.) Now suppose inductively that the procedure gives a single value to a_n . Then it gives a single value to a_{n++} , namely $a_{n++} := f_n(a_n)$. (None of the other definitions $a_{m++} := f_m(a_m)$ will redefine the value of a_{n++} , because of Axiom 2.4.) This completes the induction, and so a_n is defined for each natural number n , with a single value assigned to each a_n . \square

Note how all of the axioms had to be used here. In a system which had some sort of wrap-around, recursive definitions would not work because some elements of the sequence would constantly be redefined. For instance, in Example 2.1.5, in which $3++ = 0$, then there would be (at least) two conflicting definitions for a_0 ,

³Strictly speaking, this proposition requires one to define the notion of a function, which we shall do in the next chapter. However, this will not be circular, as the concept of a function does not require the Peano axioms. The proposition can be formalized in the language of set theory, see Exercise 3.5.12.

either c or $f_3(a_3)$). In a system which had superfluous elements such as 0.5 , the element $a_{0.5}$ would never be defined.

Recursive definitions are very powerful; for instance, we can use them to define addition and multiplication, to which we now turn.

2.2 Addition

The natural number system is very sparse right now: we have only one operation - increment - and a handful of axioms. But now we can build up more complex operations, such as addition.

The way it works is the following. To add three to five should be the same as incrementing five three times - this is one increment more than adding two to five, which is one increment more than adding one to five, which is one increment more than adding zero to five, which should just give five. So we give a recursive definition for addition as follows.

Definition 2.2.1 (Addition of natural numbers). Let m be a natural number. To add zero to m , we define $0 + m := m$. Now suppose inductively that we have defined how to add n to m . Then we can add $n++$ to m by defining $(n++) + m := (n + m)++$.

Thus $0 + m$ is m , $1 + m = (0++) + m$ is $m++$; $2 + m = (1++) + m = (m++)++$; and so forth; for instance we have $2 + 3 = (3++)++ = 4++ = 5$. From our discussion of recursion in the previous section we see that we have defined $n + m$ for every integer n (here we are specializing the previous general discussion to the setting where $a_n = n + m$ and $f_n(a_n) = a_n++$). Note that this definition is asymmetric: $3 + 5$ is incrementing 5 three times, while $5 + 3$ is incrementing 3 five times. Of course, they both yield the same value of 8. More generally, it is a fact (which we shall prove shortly) that $a + b = b + a$ for all natural numbers a, b , although this is not immediately clear from the definition.

Notice that we can prove easily, using Axioms 2.1, 2.2, and induction (Axiom 2.5), that the sum of two natural numbers is again a natural number (why?).

Right now we only have two facts about addition: that $0+m = m$, and that $(n++)+m = (n+m)++$. Remarkably, this turns out to be enough to deduce everything else we know about addition. We begin with some basic lemmas⁴.

Lemma 2.2.2. *For any natural number n , $n+0 = n$.*

Note that we cannot deduce this immediately from $0+m = m$ because we do not know yet that $a+b = b+a$.

Proof. We use induction. The base case $0+0 = 0$ follows since we know that $0+m = m$ for every natural number m , and 0 is a natural number. Now suppose inductively that $n+0 = n$. We wish to show that $(n++)+0 = n++$. But by definition of addition, $(n++)+0$ is equal to $(n+0)++$, which is equal to $n++$ since $n+0 = n$. This closes the induction. \square

Lemma 2.2.3. *For any natural numbers n and m , $n+(m++) = (n+m)++$.*

Again, we cannot deduce this yet from $(n++)+m = (n+m)++$ because we do not know yet that $a+b = b+a$.

Proof. We induct on n (keeping m fixed). We first consider the base case $n = 0$. In this case we have to prove $0+(m++) = (0+m)++$. But by definition of addition, $0+(m++) = m++$ and $0+m = m$, so both sides are equal to $m++$ and are thus equal to each other. Now we assume inductively that $n+(m++) = (n+m)++$; we now have to show that $(n++)+(m++) = ((n++)+m)++$. The left-hand side is $(n+(m++))++$ by definition of addition, which

⁴From a logical point of view, there is no difference between a lemma, proposition, theorem, or corollary - they are all claims waiting to be proved. However, we use these terms to suggest different levels of importance and difficulty. A lemma is an easily proved claim which is helpful for proving other propositions and theorems, but is usually not particularly interesting in its own right. A proposition is a statement which is interesting in its own right, while a theorem is a more important statement than a proposition which says something definitive on the subject, and often takes more effort to prove than a proposition or lemma. A corollary is a quick consequence of a proposition or theorem that was proven recently.

is equal to $((n+m)++)++$ by the inductive hypothesis. Similarly, we have $(n++)+m = (n+m)++$ by the definition of addition, and so the right-hand side is also equal to $((n+m)++)++$. Thus both sides are equal to each other, and we have closed the induction. \square

As a particular corollary of Lemma 2.2.2 and Lemma 2.2.3 we see that $n++ = n + 1$ (why?).

As promised earlier, we can now prove that $a + b = b + a$.

Proposition 2.2.4 (Addition is commutative). *For any natural numbers n and m , $n + m = m + n$.*

Proof. We shall use induction on n (keeping m fixed). First we do the base case $n = 0$, i.e., we show $0 + m = m + 0$. By the definition of addition, $0 + m = m$, while by Lemma 2.2.2, $m + 0 = m$. Thus the base case is done. Now suppose inductively that $n + m = m + n$, now we have to prove that $(n++) + m = m + (n++)$ to close the induction. By the definition of addition, $(n++) + m = (n+m)++$. By Lemma 2.2.3, $m + (n++) = (m + n)++$, but this is equal to $(n + m)++$ by the inductive hypothesis $n + m = m + n$. Thus $(n++) + m = m + (n++)$ and we have closed the induction. \square

Proposition 2.2.5 (Addition is associative). *For any natural numbers a, b, c , we have $(a + b) + c = a + (b + c)$.*

Proof. See Exercise 2.2.1. \square

Because of this associativity we can write sums such as $a + b + c$ without having to worry about which order the numbers are being added together.

Now we develop a cancellation law.

Proposition 2.2.6 (Cancellation law). *Let a, b, c be natural numbers such that $a + b = a + c$. Then we have $b = c$.*

Note that we cannot use subtraction or negative numbers yet to prove this Proposition, because we have not developed these concepts yet. In fact, this cancellation law is crucial in letting us define subtraction (and the integers) later on in these notes,

because it allows for a sort of “virtual subtraction” even before subtraction is officially defined.

Proof. We prove this by induction on a . First consider the base case $a = 0$. Then we have $0 + b = 0 + c$, which by definition of addition implies that $b = c$ as desired. Now suppose inductively that we have the cancellation law for a (so that $a + b = a + c$ implies $b = c$); we now have to prove the cancellation law for $a++$. In other words, we assume that $(a++) + b = (a++) + c$ and need to show that $b = c$. By the definition of addition, $(a++) + b = (a + b)++$ and $(a++) + c = (a + c)++$ and so we have $(a + b)++ = (a + c)++$. By Axiom 2.4, we have $a + b = a + c$. Since we already have the cancellation law for a , we thus have $b = c$ as desired. This closes the induction. \square

We now discuss how addition interacts with positivity.

Definition 2.2.7 (Positive natural numbers). A natural number n is said to be *positive* iff it is not equal to 0. (“iff” is shorthand for “if and only if”).

Proposition 2.2.8. *If a is positive and b is a natural number, then $a + b$ is positive (and hence $b + a$ is also, by Proposition 2.2.4).*

Proof. We use induction on b . If $b = 0$, then $a + b = a + 0 = a$, which is positive, so this proves the base case. Now suppose inductively that $a + b$ is positive. Then $a + (b++) = (a + b)++$, which cannot be zero by Axiom 2.3, and is hence positive. This closes the induction. \square

Corollary 2.2.9. *If a and b are natural numbers such that $a + b = 0$, then $a = 0$ and $b = 0$.*

Proof. Suppose for contradiction that $a \neq 0$ or $b \neq 0$. If $a \neq 0$ then a is positive, and hence $a + b = 0$ is positive by Proposition 2.2.8, contradiction. Similarly if $b \neq 0$ then b is positive, and again $a + b = 0$ is positive by Proposition 2.2.8, contradiction. Thus a and b must both be zero. \square

Lemma 2.2.10. *Let a be a positive number. Then there exists exactly one natural number b such that $b++ = a$.*

Proof. See Exercise 2.2.2. □

Once we have a notion of addition, we can begin defining a notion of *order*.

Definition 2.2.11 (Ordering of the natural numbers). Let n and m be natural numbers. We say that n is *greater than or equal to* m , and write $n \geq m$ or $m \leq n$, iff we have $n = m + a$ for some natural number a . We say that n is *strictly greater than* m , and write $n > m$ or $m < n$, iff $n \geq m$ and $n \neq m$.

Thus for instance $8 > 5$, because $8 = 5 + 3$ and $8 \neq 5$. Also note that $n++ > n$ for any n ; thus there is no largest natural number n , because the next number $n++$ is always larger still.

Proposition 2.2.12 (Basic properties of order for natural numbers). *Let a, b, c be natural numbers. Then*

- (Order is reflexive) $a \geq a$.
- (Order is transitive) If $a \geq b$ and $b \geq c$, then $a \geq c$.
- (Order is anti-symmetric) If $a \geq b$ and $b \geq a$, then $a = b$.
- (Addition preserves order) $a \geq b$ if and only if $a + c \geq b + c$.
- $a < b$ if and only if $a++ \leq b$.
- $a < b$ if and only if $b = a + d$ for some positive number d .

Proof. See Exercise 2.2.3. □

Proposition 2.2.13 (Trichotomy of order for natural numbers). *Let a and b be natural numbers. Then exactly one of the following statements is true: $a < b$, $a = b$, or $a > b$.*

Proof. This is only a sketch of the proof; the gaps will be filled in Exercise 2.2.4.

First we show that we cannot have more than one of the statements $a < b$, $a = b$, $a > b$ holding at the same time. If $a < b$ then $a \neq b$ by definition, and if $a > b$ then $a \neq b$ by definition. If $a > b$ and $a < b$ then by Proposition 2.2.12 we have $a = b$, a contradiction. Thus no more than one of the statements is true.

Now we show that at least one of the statements is true. We keep b fixed and induct on a . When $a = 0$ we have $0 \leq b$ for all b (why?), so we have either $0 = b$ or $0 < b$, which proves the base case. Now suppose we have proven the Proposition for a , and now we prove the proposition for $a++$. From the trichotomy for a , there are three cases: $a < b$, $a = b$, and $a > b$. If $a > b$, then $a++ > b$ (why?). If $a = b$, then $a++ > b$ (why?). Now suppose that $a < b$. Then by Proposition 2.2.12, we have $a++ \leq b$. Thus either $a++ = b$ or $a++ < b$, and in either case we are done. This closes the induction. \square

The properties of order allow one to obtain a stronger version of the principle of induction:

Proposition 2.2.14 (Strong principle of induction). *Let m_0 be a natural number, and Let $P(m)$ be a property pertaining to an arbitrary natural number m . Suppose that for each $m \geq m_0$, we have the following implication: if $P(m')$ is true for all natural numbers $m_0 \leq m' < m$, then $P(m)$ is also true. (In particular, this means that $P(m_0)$ is true, since in this case the hypothesis is vacuous.) Then we can conclude that $P(m)$ is true for all natural numbers $m \geq m_0$.*

Remark 2.2.15. In applications we usually use this principle with $m_0 = 0$ or $m_0 = 1$.

Proof. See Exercise 2.2.5. \square

Exercise 2.2.1. Prove Proposition 2.2.5. (Hint: fix two of the variables and induct on the third.)

Exercise 2.2.2. Prove Lemma 2.2.10. (Hint: use induction.)

Exercise 2.2.3. Prove Proposition 2.2.12. (Hint: you will need many of the preceding propositions, corollaries, and lemmas.)

Exercise 2.2.4. Justify the three statements marked (why?) in the proof of Proposition 2.2.13.

Exercise 2.2.5. Prove Proposition 2.2.14. (Hint: define $Q(n)$ to be the property that $P(m)$ is true for all $m_0 \leq m < n$; note that $Q(n)$ is vacuously true when $n < m_0$.)

Exercise 2.2.6. Let n be a natural number, and let $P(m)$ be a property pertaining to the natural numbers such that whenever $P(m++)$ is true, then $P(m)$ is true. Suppose that $P(n)$ is also true. Prove that $P(m)$ is true for all natural numbers $m \leq n$; this is known as the *principle of backwards induction*. (Hint: Apply induction to the variable n .)

2.3 Multiplication

In the previous section we have proven all the basic facts that we know to be true about addition and order. To save space and to avoid belaboring the obvious, we will now allow ourselves to use all the rules of algebra concerning addition and order that we are familiar with, without further comment. Thus for instance we may write things like $a + b + c = c + b + a$ without supplying any further justification. Now we introduce multiplication. Just as addition is the iterated increment operation, multiplication is iterated addition:

Definition 2.3.1 (Multiplication of natural numbers). Let m be a natural number. To multiply zero to m , we define $0 \times m := 0$. Now suppose inductively that we have defined how to multiply n to m . Then we can multiply $n++$ to m by defining $(n++) \times m := (n \times m) + m$.

Thus for instance $0 \times m = 0$, $1 \times m = 0 + m$, $2 \times m = 0 + m + m$, etc. By induction one can easily verify that the product of two natural numbers is a natural number.

Lemma 2.3.2 (Multiplication is commutative). *Let n, m be natural numbers. Then $n \times m = m \times n$.*

Proof. See Exercise 2.3.1. □

We will now abbreviate $n \times m$ as nm , and use the usual convention that multiplication takes precedence over addition, thus for instance $ab + c$ means $(a \times b) + c$, not $a \times (b + c)$. (We will also use the usual notational conventions of precedence for the other arithmetic operations when they are defined later, to save on using parentheses all the time.)

Lemma 2.3.3 (Natural numbers have no zero divisors). *Let n, m be natural numbers. Then $n \times m = 0$ if and only if at least one of n, m is equal to zero. In particular, if n and m are both positive, then nm is also positive.*

Proof. See Exercise 2.3.2. □

Proposition 2.3.4 (Distributive law). *For any natural numbers a, b, c , we have $a(b + c) = ab + ac$ and $(b + c)a = ba + ca$.*

Proof. Since multiplication is commutative we only need to show the first identity $a(b + c) = ab + ac$. We keep a and b fixed, and use induction on c . Let's prove the base case $c = 0$, i.e., $a(b + 0) = ab + a0$. The left-hand side is ab , while the right-hand side is $ab + 0 = ab$, so we are done with the base case. Now let us suppose inductively that $a(b + c) = ab + ac$, and let us prove that $a(b + (c++)) = ab + a(c++)$. The left-hand side is $a((b + c)++) = a(b + c) + a$, while the right-hand side is $ab + ac + a = a(b + c) + a$ by the induction hypothesis, and so we can close the induction. □

Proposition 2.3.5 (Multiplication is associative). *For any natural numbers a, b, c , we have $(a \times b) \times c = a \times (b \times c)$.*

Proof. See Exercise 2.3.3. □

Proposition 2.3.6 (Multiplication preserves order). *If a, b are natural numbers such that $a < b$, and c is positive, then $ac < bc$.*

Proof. Since $a < b$, we have $b = a + d$ for some positive d . Multiplying by c and using the distributive law we obtain $bc = ac + dc$. Since d is positive, and c is non-zero (hence positive), dc is positive, and hence $ac < bc$ as desired. \square

Corollary 2.3.7 (Cancellation law). *Let a, b, c be natural numbers such that $ac = bc$ and c is non-zero. Then $a = b$.*

Remark 2.3.8. Just as Proposition 2.2.6 will allow for a “virtual subtraction” which will eventually let us define genuine subtraction, this corollary provides a “virtual division” which will be needed to define genuine division later on.

Proof. By the trichotomy of order (Proposition 2.2.13), we have three cases: $a < b$, $a = b$, $a > b$. Suppose first that $a < b$, then by Proposition 2.3.6 we have $ac < bc$, a contradiction. We can obtain a similar contradiction when $a > b$. Thus the only possibility is that $a = b$, as desired. \square

With these propositions it is easy to deduce all the familiar rules of algebra involving addition and multiplication, see for instance Exercise 2.3.4.

Now that we have the familiar operations of addition and multiplication, the more primitive notion of increment will begin to fall by the wayside, and we will see it rarely from now on. In any event we can always use addition to describe incrementation, since $n++ = n + 1$.

Proposition 2.3.9 (Euclidean algorithm). *Let n be a natural number, and let q be a positive number. Then there exist natural numbers m, r such that $0 \leq r < q$ and $n = mq + r$.*

Remark 2.3.10. In other words, we can divide a natural number n by a positive number q to obtain a quotient m (which is another natural number) and a remainder r (which is less than q). This algorithm marks the beginning of *number theory*, which is a beautiful and important subject but one which is beyond the scope of this text.

Proof. See Exercise 2.3.5. □

Just like one uses the increment operation to recursively define addition, and addition to recursively define multiplication, one can use multiplication to recursively define *exponentiation*:

Definition 2.3.11 (Exponentiation for natural numbers). Let m be a natural number. To raise m to the power 0, we define $m^0 := 1$. Now suppose recursively that m^n has been defined for some natural number n , then we define $m^{n++} := m^n \times m$.

Examples 2.3.12. Thus for instance $x^1 = x^0 \times x = 1 \times x = x$; $x^2 = x^1 \times x = x \times x$; $x^3 = x^2 \times x = x \times x \times x$; and so forth. By induction we see that this recursive definition defines x^n for all natural numbers n .

We will not develop the theory of exponentiation too deeply here, but instead wait until after we have defined the integers and rational numbers; see in particular Proposition 4.3.10.

Exercise 2.3.1. Prove Lemma 2.3.2. (Hint: modify the proofs of Lemmas 2.2.2, 2.2.3 and Proposition 2.2.4).

Exercise 2.3.2. Prove Lemma 2.3.3. (Hint: prove the second statement first).

Exercise 2.3.3. Prove Proposition 2.3.5. (Hint: modify the proof of Proposition 2.2.5 and use the distributive law).

Exercise 2.3.4. Prove the identity $(a + b)^2 = a^2 + 2ab + b^2$ for all natural numbers a, b .

Exercise 2.3.5. Prove Proposition 2.3.9. (Hint: Fix q and induct on n).

Chapter 3

Set theory

Modern analysis, like most of modern mathematics, is concerned with numbers, sets, and geometry. We have already introduced one type of number system, the natural numbers. We will introduce the other number systems shortly, but for now we pause to introduce the concepts and notation of set theory, as they will be used increasingly heavily in later chapters. (We will not pursue a rigorous description of Euclidean geometry in this text, preferring instead to describe that geometry in terms of the real number system by means of the Cartesian co-ordinate system.)

While set theory is not the main focus of this text, almost every other branch of mathematics relies on set theory as part of its foundation, so it is important to get at least some grounding in set theory before doing other advanced areas of mathematics. In this chapter we present the more elementary aspects of axiomatic set theory, leaving more advanced topics such as a discussion of infinite sets and the axiom of choice to Chapter 8. A full treatment of the finer subtleties of set theory (of which there are many!) is unfortunately well beyond the scope of this text.

3.1 Fundamentals

In this section we shall set out some axioms for sets, just as we did for the natural numbers. For pedagogical reasons, we will use a somewhat overcomplete list of axioms for set theory, in the sense

that some of the axioms can be used to deduce others, but there is no real harm in doing this. We begin with an informal description of what sets should be.

Definition 3.1.1. (Informal) We define a *set* A to be any unordered collection of objects, e.g., $\{3, 8, 5, 2\}$ is a set. If x is an object, we say that x is an *element of* A or $x \in A$ if x lies in the collection; otherwise we say that $x \notin S$. For instance, $3 \in \{1, 2, 3, 4, 5\}$ but $7 \notin \{1, 2, 3, 4, 5\}$.

This definition is intuitive enough, but it doesn't answer a number of questions, such as which collections of objects are considered to be sets, which sets are equal to other sets, and how one defines operations on sets (e.g., unions, intersections, etc.). Also, we have no axioms yet on what sets do, or what their elements do. Obtaining these axioms and defining these operations will be the purpose of the remainder of this section.

We first clarify one point: we consider sets themselves to be a type of object.

Axiom 3.1 (Sets are objects). *If A is a set, then A is also an object. In particular, given two sets A and B , it is meaningful to ask whether A is also an element of B .*

Example 3.1.2. (Informal) The set $\{3, \{3, 4\}, 4\}$ is a set of three distinct elements, one of which happens to itself be a set of two elements. See Example 3.1.10 for a more formal version of this example. However, not all objects are sets; for instance, we typically do not consider a natural number such as 3 to be a set. (The more accurate statement is that natural numbers can be the *cardinalities* of sets, rather than necessarily being sets themselves. See Section 3.6.)

Remark 3.1.3. There is a special case of set theory, called “pure set theory”, in which *all* objects are sets; for instance the number 0 might be identified with the empty set $\emptyset = \{\}$, the number 1 might be identified with $\{0\} = \{\{\}\}$, the number 2 might be identified with $\{0, 1\} = \{\{\}, \{\{\}\}\}$, and so forth. From a logical point of

view, pure set theory is a simpler theory, since one only has to deal with sets and not with objects; however, from a conceptual point of view it is often easier to deal with impure set theories in which some objects are not considered to be sets. The two types of theories are more or less equivalent for the purposes of doing mathematics, and so we shall take an agnostic position as to whether all objects are sets or not.

To summarize so far, among all the objects studied in mathematics, some of the objects happen to be sets; and if x is an object and A is a set, then either $x \in A$ is true or $x \in A$ is false. (If A is not a set, we leave the statement $x \in A$ undefined; for instance, we consider the statement $3 \in 4$ to neither be true or false, but simply meaningless, since 4 is not a set.)

Next, we define the notion of equality: when are two sets considered to be equal? We do not consider the order of the elements inside a set to be important; thus we think of $\{3, 8, 5, 2\}$ and $\{2, 3, 5, 8\}$ as the same set. On the other hand, $\{3, 8, 5, 2\}$ and $\{3, 8, 5, 2, 1\}$ are different sets, because the latter set contains an element (1) that the former one does not. For similar reasons $\{3, 8, 5, 2\}$ and $\{3, 8, 5\}$ are different sets. We formalize this as a definition:

Definition 3.1.4 (Equality of sets). Two sets A and B are *equal*, $A = B$, iff every element of A is an element of B and vice versa. To put it another way, $A = B$ if and only if every element x of A belongs also to B , and every element y of B belongs also to A .

Example 3.1.5. Thus, for instance, $\{1, 2, 3, 4, 5\}$ and $\{3, 4, 2, 1, 5\}$ are the same set, since they contain exactly the same elements. (The set $\{3, 3, 1, 5, 2, 4, 2\}$ is also equal to $\{1, 2, 3, 4, 5\}$; the additional repetition of 3 and 2 is redundant as it does not further change the status of 2 and 3 being elements of the set.)

One can easily verify that this notion of equality is reflexive, symmetric, and transitive (Exercise 3.1.1). Observe that if $x \in A$ and $A = B$, then $x \in B$, by Definition 3.1.4. Thus the “is an element of” relation \in obeys the axiom of substitution (see Section

12.7). Because of this, any new operation we define on sets will also obey the axiom of substitution, as long as we can define that operation purely in terms of the relation \in . This is for instance the case for the remaining definitions in this section. (On the other hand, we cannot use the notion of the “first” or “last” element in a set in a well-defined manner, because this would not respect the axiom of substitution; for instance the sets $\{1, 2, 3, 4, 5\}$ and $\{3, 4, 2, 1, 5\}$ are the same set, but have different first elements.)

Next, we turn to the issue of exactly which objects are sets and which objects are not. The situation is analogous to how we defined the natural numbers in the previous chapter; we started with a single natural number, 0, and started building more numbers out of 0 using the increment operation. We will try something similar here, starting with a single set, the *empty set*, and building more sets out of the empty set by various operations. We begin by postulating the existence of the empty set.

Axiom 3.2 (Empty set). *There exists a set \emptyset , known as the empty set, which contains no elements, i.e., for every object x we have $x \notin \emptyset$.*

The empty set is also denoted $\{\}$. Note that there can only be one empty set; if there were two sets \emptyset and \emptyset' which were both empty, then by Definition 3.1.4 they would be equal to each other (why?).

If a set is not equal to the empty set, we call it *non-empty*. The following statement is very simple, but worth stating nevertheless:

Lemma 3.1.6 (Single choice). *Let A be a non-empty set. Then there exists an object x such that $x \in A$.*

Proof. We prove by contradiction. Suppose there does not exist any object x such that $x \in A$. Then for all objects x , $x \notin A$. Also, by Axiom 3.2 we have $x \notin \emptyset$. Thus $x \in A \iff x \in \emptyset$ (both statements are equally false), and so $A = \emptyset$ by Definition 3.1.4. \square

Remark 3.1.7. The above Lemma asserts that given any non-empty set A , we are allowed to “choose” an element x of A which

demonstrates this non-emptiness. Later on (in Lemma 3.5.12) we will show that given any finite number of non-empty sets, say A_1, \dots, A_n , it is possible to choose one element x_1, \dots, x_n from each set A_1, \dots, A_n ; this is known as “finite choice”. However, in order to choose elements from an infinite number of sets, we need an additional axiom, the *axiom of choice*, which we will discuss in Section 8.4.

Remark 3.1.8. Note that the empty set is *not* the same thing as the natural number 0. One is a set; the other is a number. However, it is true that the *cardinality* of the empty set is 0; see Section 3.6.

If Axiom 3.2 was the only axiom that set theory had, then set theory could be quite boring, as there might be just a single set in existence, the empty set. We now make some further axioms to enrich the class of sets we can make.

Axiom 3.3 (Singleton sets and pair sets). *If a is an object, then there exists a set $\{a\}$ whose only element is a , i.e., for every object y , $y \in \{a\}$ if and only if $y = a$; we refer to $\{a\}$ as the singleton set whose element is a . Furthermore, if a and b are objects, then there exists a set $\{a, b\}$ whose only elements are a and b ; i.e., for every object y , $y \in \{a, b\}$ if and only if $y = a$ or $y = b$; we refer to this set as the pair set formed by a and b .*

Remarks 3.1.9. Just as there is only one empty set, there is only one singleton set for each object a , thanks to Definition 3.1.4 (why?). Similarly, given any two objects a and b , there is only one pair set formed by a and b . Also, Definition 3.1.4 also ensures that $\{a, b\} = \{b, a\}$ (why?) and $\{a, a\} = \{a\}$ (why?). Thus the singleton set axiom is in fact redundant, being a consequence of the pair set axiom. Conversely, the pair set axiom will follow from the singleton set axiom and the pairwise union axiom below (see Lemma 3.1.12). One may wonder why we don’t go further and create triplet axioms, quadruplet axioms, etc.; however there will be no need for this once we introduce the pairwise union axiom below.

Examples 3.1.10. Since \emptyset is a set (and hence an object), so is singleton set $\{\emptyset\}$, i.e., the set whose only element is \emptyset , is a set (and it is *not* the same set as \emptyset , $\{\emptyset\} \neq \emptyset$ (why?). So is the singleton set $\{\{\emptyset\}\}$ and the pair set $\{\emptyset, \{\emptyset\}\}$ is also a set. These three sets are all unequal to each other (Exercise 3.1.2).

As the above examples show, we now can create quite a few sets; however, the sets we make are still fairly small (each set that we can build consists of no more than two elements, so far). The next axiom allows us to build somewhat larger sets than before.

Axiom 3.4 (Pairwise union). *Given any two sets A, B , there exists a set $A \cup B$, called the union $A \cup B$ of A and B , whose elements consists of all the elements which belong to A or B or both. In other words, for any object x ,*

$$x \in A \cup B \iff (x \in A \text{ or } x \in B).$$

(Recall that “or” refers by default in mathematics to inclusive or: “ X or Y is true” means that “either X is true, or Y is true, or both are true”.)

Example 3.1.11. The set $\{1, 2\} \cup \{2, 3\}$ consists of those elements which either lie on $\{1, 2\}$ or in $\{2, 3\}$ or in both, or in other words the elements of this set are simply 1, 2, and 3. Because of this, we denote this set as $\{1, 2\} \cup \{2, 3\} = \{1, 2, 3\}$.

Note that if A, B, A' are sets, and A is equal to A' , then $A \cup B$ is equal to $A' \cup B$ (why? One needs to use Axiom 3.4 and Definition 3.1.4). Similarly if B' is a set which is equal to B , then $A \cup B$ is equal to $A \cup B'$. Thus the operation of union obeys the axiom of substitution, and is thus well-defined on sets.

We now give some basic properties of unions.

Lemma 3.1.12. *If a and b are objects, then $\{a, b\} = \{a\} \cup \{b\}$. If A, B, C are sets, then the union operation is commutative (i.e., $A \cup B = B \cup A$) and associative (i.e., $(A \cup B) \cup C = A \cup (B \cup C)$). Also, we have $A \cup A = A \cup \emptyset = \emptyset \cup A = A$.*

Proof. We prove just the associativity identity $(A \cup B) \cup C = A \cup (B \cup C)$, and leave the remaining claims to Exercise 3.1.3. By Definition 3.1.4, we need to show that every element x of $(A \cup B) \cup C$ is an element of $A \cup (B \cup C)$, and vice versa. So suppose first that x is an element of $(A \cup B) \cup C$. By Axiom 3.4, this means that at least one of $x \in A \cup B$ or $x \in C$ is true, so we divide into cases. If $x \in C$, then by Axiom 3.4 again $x \in B \cup C$, and so by Axiom 3.4 again $x \in A \cup (B \cup C)$. Now suppose instead $x \in A \cup B$, then by Axiom 3.4 again $x \in A$ or $x \in B$. If $x \in A$ then $x \in A \cup (B \cup C)$ by Axiom 3.4, while if $x \in B$ then by two applications of Axiom 3.4 we have $x \in B \cup C$ and hence $x \in A \cup (B \cup C)$. Thus in all cases we see that every element of $(A \cup B) \cup C$ lies in $A \cup (B \cup C)$. A similar argument shows that every element of $A \cup (B \cup C)$ lies in $(A \cup B) \cup C$, and so $(A \cup B) \cup C = A \cup (B \cup C)$ as desired. \square

Because of the above lemma, we do not need to use parentheses to denote multiple unions, thus for instance we can write $A \cup B \cup C$ instead of $(A \cup B) \cup C$ or $A \cup (B \cup C)$. Similarly for unions of four sets, $A \cup B \cup C \cup D$, etc.

Remark 3.1.13. Note that while the operation of union has some similarities with addition, the two operations are *not* identical. For instance, $\{2\} \cup \{3\} = \{2, 3\}$ and $2 + 3 = 5$, whereas $\{2\} + \{3\}$ is meaningless (addition pertains to numbers, not sets) and $2 \cup 3$ is also meaningless (union pertains to sets, not numbers).

This axiom allows us to define triplet sets, quadruplet sets, and so forth: if a, b, c are three objects, we define $\{a, b, c\} := \{a\} \cup \{b\} \cup \{c\}$; if a, b, c, d are four objects, then we define $\{a, b, c, d\} := \{a\} \cup \{b\} \cup \{c\} \cup \{d\}$, and so forth. On the other hand, we are not yet in a position to define sets consisting of n objects for any given natural number n ; this would require iterating the above construction “ n times”, but the concept of n -fold iteration has not yet been rigorously defined. For similar reasons, we cannot yet define sets consisting of infinitely many objects, because that would require iterating the axiom of pairwise union infinitely often, and it is not clear at this stage that one can do this rigorously. Later on,

we will introduce other axioms of set theory which allow one to construct arbitrarily large, and even infinite, sets.

Clearly, some sets seem to be larger than others. One way to formalize this concept is through the notion of a *subset*.

Definition 3.1.14 (Subsets). Let A, B be sets. We say that A is a *subset* of B , denoted $A \subseteq B$, iff every element of A is also an element of B , i.e.

$$\text{For any object } x, \quad x \in A \implies x \in B.$$

We say that A is a *proper subset* of B , denoted $A \subsetneq B$, if $A \subseteq B$ and $A \neq B$.

Remark 3.1.15. Note that because these definitions involve only the notions of equality and the “is an element of” relation, both of which already obey the axiom of substitution, the notion of subset also automatically obeys the axiom of substitution. Thus for instance if $A \subseteq B$ and $A = A'$, then $A' \subseteq B$.

Examples 3.1.16. We have $\{1, 2, 4\} \subseteq \{1, 2, 3, 4, 5\}$, because every element of $\{1, 2, 4\}$ is also an element of $\{1, 2, 3, 4, 5\}$. In fact we also have $\{1, 2, 4\} \subsetneq \{1, 2, 3, 4, 5\}$, since the two sets $\{1, 2, 4\}$ and $\{1, 2, 3, 4, 5\}$ are not equal. Given any set A , we always have $A \subseteq A$ (why?) and $\emptyset \subseteq A$ (why?).

The notion of subset in set theory is similar to the notion of “less than or equal to” for numbers, as the following Proposition demonstrates (for a more precise statement, see Definition 8.5.1):

Proposition 3.1.17 (Sets are partially ordered by set inclusion). *Let A, B, C be sets. If $A \subseteq B$ and $B \subseteq C$ then $A \subseteq C$. If instead $A \subseteq B$ and $B \subseteq A$, then $A = B$. Finally, if $A \subsetneq B$ and $B \subsetneq C$ then $A \subsetneq C$.*

Proof. We shall just prove the first claim. Suppose that $A \subseteq B$ and $B \subseteq C$. To prove that $A \subseteq C$, we have to prove that every element of A is an element of C . So, let us pick an arbitrary element x of A . Then, since $A \subseteq B$, x must then be an element of B . But then since $B \subseteq C$, x is an element of C . Thus every element of A is indeed an element of C , as claimed. \square

Remark 3.1.18. There is a relationship between subsets and unions; see for instance Exercise 3.1.7.

Remark 3.1.19. There is one important difference between the subset relation \subsetneq and the less than relation $<$. Given any two distinct natural numbers n, m , we know that one of them is smaller than the other (Proposition 2.2.13); however, given two distinct sets, it is not in general true that one of them is a subset of the other. For instance, take $A := \{2n : n \in \mathbf{N}\}$ to be the set of even natural numbers, and $B := \{2n + 1 : n \in \mathbf{N}\}$ to be the set of odd natural numbers. Then neither set is a subset of the other. This is why we say that sets are only *partially ordered*, whereas the natural numbers are *totally ordered* (see Definitions 8.5.1, 8.5.3).

Remark 3.1.20. We should also caution that the subset relation \subseteq is not the same as the element relation \in . The number 2 is an element of $\{1, 2, 3\}$, thus $2 \in \{1, 2, 3\}$, but is not a subset of $\{1, 2, 3\}$, $2 \not\subseteq \{1, 2, 3\}$; indeed, 2 is not even a set. Conversely, while $\{2\}$ is a subset of $\{1, 2, 3\}$, $\{2\} \subseteq \{1, 2, 3\}$, it is not an element, $\{2\} \notin \{1, 2, 3\}$. The point is that the number 2 and the set $\{2\}$ are distinct objects. It is important to distinguish sets from their elements, as they can have different properties. For instance, it is possible to have an infinite set consisting of finite numbers (the set \mathbf{N} of natural numbers is one such example), and it is also possible to have a finite set consisting of infinite objects (consider for instance the set $\{\mathbf{N}, \mathbf{Z}, \mathbf{Q}, \mathbf{R}\}$, which has four elements, all of which are infinite).

We now give an axiom which easily allows us to create subsets out of larger sets.

Axiom 3.5 (Axiom of specification). *Let A be a set, and for each $x \in A$, let $P(x)$ be a property pertaining to x (i.e., $P(x)$ is either a true statement or a false statement). Then there exists a set, called $\{x \in A : P(x) \text{ is true}\}$ (or simply $\{x \in A : P(x)\}$ for short), whose elements are precisely the elements x in A for which $P(x)$ is true. In other words, for any object y ,*

$$y \in \{x \in A : P(x) \text{ is true}\} \iff (y \in A \text{ and } P(y) \text{ is true.})$$

This axiom is also known as the *axiom of separation*. Note that $\{x \in A : P(x) \text{ is true}\}$ is always a subset of A (why?), though it could be as large as A or as small as the empty set. One can verify that the axiom of substitution works for specification, thus if $A = A'$ then $\{x \in A : P(x)\} = \{x \in A' : P(x)\}$ (why?).

Example 3.1.21. Let $S := \{1, 2, 3, 4, 5\}$. Then the set $\{n \in S : n < 4\}$ is the set of those elements n in S for which $n < 4$ is true, i.e., $\{n \in S : n < 4\} = \{1, 2, 3\}$. Similarly, the set $\{n \in S : n < 7\}$ is the same as S itself, while $\{n \in S : n < 1\}$ is the empty set.

We sometimes write $\{x \in A \mid P(x)\}$ instead of $\{x \in A : P(x)\}$; this is useful when we are using the colon “:” to denote something else, for instance to denote the range and domain of a function $f : X \rightarrow Y$.

We can use this axiom of specification to define some further operations on sets, namely intersections and difference sets.

Definition 3.1.22 (Intersections). The *intersection* $S_1 \cap S_2$ of two sets is defined to be the set

$$S_1 \cap S_2 := \{x \in S_1 : x \in S_2\}.$$

In other words, $S_1 \cap S_2$ consists of all the elements which belong to both S_1 and S_2 . Thus, for all objects x ,

$$x \in S_1 \cap S_2 \iff x \in S_1 \text{ and } x \in S_2.$$

Remark 3.1.23. Note that this definition is well-defined (i.e., it obeys the axiom of substitution, see Section 12.7) because it is defined in terms of more primitive operations which were already known to obey the axiom of substitution. Similar remarks apply to future definitions in this chapter and will usually not be mentioned explicitly again.

Examples 3.1.24. We have $\{1, 2, 4\} \cap \{2, 3, 4\} = \{2, 4\}$, $\{1, 2\} \cap \{3, 4\} = \emptyset$, $\{2, 3\} \cup \emptyset = \{2, 3\}$, and $\{2, 3\} \cap \emptyset = \emptyset$.

Remark 3.1.25. By the way, one should be careful with the English word “and”: rather confusingly, it can mean either union or intersection, depending on context. For instance, if one talks about a set of “boys and girls”, one means the *union* of a set of boys with a set of girls, but if one talks about the set of people who are single and male, then one means the *intersection* of the set of single people with the set of male people. (Can you work out the rule of grammar that determines when “and” means union and when “and” means intersection?) Another problem is that “and” is also used in English to denote addition, thus for instance one could say that “2 and 3 is 5”, while also saying that “the elements of $\{2\}$ and the elements of $\{3\}$ form the set $\{2, 3\}$ ” and “the elements in $\{2\}$ and $\{3\}$ form the set \emptyset ”. This can certainly get confusing! One reason we resort to mathematical symbols instead of English words such as “and” is that mathematical symbols always have a precise and unambiguous meaning, whereas one must often look very carefully at the context in order to work out what an English word means.

Two sets A, B are said to be *disjoint* if $A \cap B = \emptyset$. Note that this is not the same concept as being *distinct*, $A \neq B$. For instance, the sets $\{1, 2, 3\}$ and $\{2, 3, 4\}$ are distinct (there are elements of one set which are not elements of the other) but not disjoint (because their intersection is non-empty). Meanwhile, the sets \emptyset and \emptyset are disjoint but not distinct (why?).

Definition 3.1.26 (Difference sets). Given two sets A and B , we define the set $A - B$ or $A \setminus B$ to be the set A with any elements of B removed:

$$A - B := \{x \in A : x \notin B\};$$

for instance, $\{1, 2, 3, 4\} \setminus \{2, 4, 6\} = \{1, 3\}$. In many cases B will be a subset of A , but not necessarily.

We now give some basic properties of unions, intersections, and difference sets.

Proposition 3.1.27 (Sets form a boolean algebra). *Let A, B, C be sets, and let X be a set containing A, B, C as subsets.*

- (*Minimal element*) We have $A \cup \emptyset = A$ and $A \cap \emptyset = \emptyset$.
- (*Maximal element*) We have $A \cup X = X$ and $A \cap X = A$.
- (*Identity*) We have $A \cap A = A$ and $A \cup A = A$.
- (*Commutativity*) We have $A \cup B = B \cup A$ and $A \cap B = B \cap A$.
- (*Associativity*) We have $(A \cup B) \cup C = A \cup (B \cup C)$ and $(A \cap B) \cap C = A \cap (B \cap C)$.
- (*Distributivity*) We have $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$ and $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$.
- (*Partition*) We have $A \cup (X \setminus A) = X$ and $A \cap (X \setminus A) = \emptyset$.
- (*De Morgan laws*) We have $X \setminus (A \cup B) = (X \setminus A) \cap (X \setminus B)$ and $X \setminus (A \cap B) = (X \setminus A) \cup (X \setminus B)$.

Remark 3.1.28. The de Morgan laws are named after the logician Augustus De Morgan (1806–1871), who identified them as one of the basic laws of set theory.

Proof. We will prove just one law, that $A \cup B = B \cup A$, and leave the remainder as an exercise (Exercise 3.1.6). We have to show that every element of $A \cup B$ is an element of $B \cup A$ and vice versa. But if $x \in A \cup B$, then by Axiom 3.4 we know that x belongs to A or to B or both. In either case it is clear that x belongs to B or A or to both, hence $x \in B \cup A$. Similarly if $x \in B \cup A$ then $x \in A \cup B$, and so the two sets are equal. \square

Remark 3.1.29. The reader may observe a certain symmetry in the above laws between \cup and \cap , and between X and \emptyset . This is an example of *duality* - two distinct properties or objects being dual to each other. In this case, the duality is manifested by the complementation relation $A \mapsto X \setminus A$; the de Morgan laws assert that this relation converts unions into intersections and vice versa. (It also interchanges X and the empty set). The above laws are collectively known as the *laws of Boolean algebra*, after the mathematician George Boole (1815–1864), and are also applicable

to a number of other objects other than sets; it plays a particularly important role in logic.

We have now accumulated a number of axioms and results about sets, but there are still many things we are not able to do yet. One of the basic things we wish to do with a set is take each of the objects of that set, and somehow transform each such object into a new object; for instance we may wish to start with a set of numbers, say $\{3, 5, 9\}$, and increment each one, creating a new set $\{4, 6, 10\}$. This is not something we can do directly using only the axioms we already have, so we need a new axiom:

Axiom 3.6 (Replacement). *Let A be a set. For any object $x \in A$, and any object y , suppose we have a statement $P(x, y)$ pertaining to x and y , such that for each $x \in A$ there is at most one y for which $P(x, y)$ is true. Then there exists a set $\{y : P(x, y) \text{ is true for some } x \in A\}$, such that for any object z ,*

$$\begin{aligned} z \in \{y : P(x, y) \text{ is true for some } x \in A\} \\ \iff P(x, z) \text{ is true for some } x \in A. \end{aligned}$$

Example 3.1.30. Let $A := \{3, 5, 9\}$, and let $P(x, y)$ be the statement $y = x++$, i.e., y is the successor of x . Observe that for every $x \in A$, there is exactly one y for which $P(x, y)$ is true - specifically, the successor of x . Thus the above axiom asserts that the set $\{y : y = x++ \text{ for some } x \in \{3, 5, 9\}\}$ exists; in this case, it is clearly the same set as $\{4, 6, 10\}$ (why?).

Example 3.1.31. Let $A = \{3, 5, 9\}$, and let $P(x, y)$ be the statement $y = 1$. Then again for every $x \in A$, there is exactly one y for which $P(x, y)$ is true - specifically, the number 1. In this case $\{y : y = 1 \text{ for some } x \in \{3, 5, 9\}\}$ is just the singleton set $\{1\}$; we have replaced each element 3, 5, 9 of the original set A by the same object, namely 1. Thus this rather silly example shows that the set obtained by the above axiom can be “smaller” than the original set.

We often abbreviate a set of the form $\{y : y = f(x) \text{ for some } x \in A\}$ as $\{f(x) : x \in A\}$ or $\{f(x) \mid x \in A\}$. Thus for instance, if

$A = \{3, 5, 9\}$, then $\{x++ : x \in A\}$ is the set $\{4, 6, 10\}$. We can of course combine the axiom of replacement with the axiom of specification, thus for instance we can create sets such as $\{f(x) : x \in A; P(x) \text{ is true}\}$ by starting with the set A , using the axiom of specification to create the set $\{x \in A : P(x) \text{ is true}\}$, and then applying the axiom of replacement to create $\{f(x) : x \in A; P(x) \text{ is true}\}$. Thus for instance $\{n++ : n \in \{3, 5, 9\}; n < 6\} = \{4, 6\}$.

In many of our examples we have implicitly assumed that natural numbers are in fact objects. Let us formalize this as follows.

Axiom 3.7 (Infinity). *There exists a set \mathbf{N} , whose elements are called natural numbers, as well as an object 0 in \mathbf{N} , and an object $n++$ assigned to every natural number $n \in \mathbf{N}$, such that the Peano axioms (Axiom 2.1 - 2.5) hold.*

This is the more formal version of Assumption 2.6. It is called the axiom of infinity because it introduces the most basic example of an infinite set, namely the set of natural numbers \mathbf{N} . (We will formalize what finite and infinite mean in Section 3.6). From the axiom of infinity we see that numbers such as 3, 5, 7, etc. are indeed objects in set theory, and so (from the pair set axiom and pairwise union axiom) we can indeed legitimately construct sets such as $\{3, 5, 9\}$ as we have been doing in our examples.

One has to keep the concept of a set distinct from the elements of that set; for instance, the set $\{n + 3 : n \in \mathbf{N}, 0 \leq n \leq 5\}$ is not the same thing as the expression or function $n + 3$. We emphasize this with an example:

Example 3.1.32. (Informal) This example requires the notion of subtraction, which has not yet been formally introduced. The following two sets are equal,

$$\{n + 3 : n \in \mathbf{N}, 0 \leq n \leq 5\} = \{8 - n : n \in \mathbf{N}, 0 \leq n \leq 5\}, \quad (3.1)$$

(see below), even though the expressions $n + 3$ and $8 - n$ are never equal to each other for any natural number n . Thus, it is a good idea to remember to put those curly braces $\{\}$ in when you talk

about sets, lest you accidentally confuse a set with its elements. One reason for this counter-intuitive situation is that the letter n is being used in two different ways on the two different sides of (3.1). To clarify the situation, let us rewrite the set $\{8 - n : n \in \mathbf{N}, 0 \leq n \leq 5\}$ by replacing the letter n by the letter m , thus giving $\{8 - m : m \in \mathbf{N}, 0 \leq m \leq 5\}$. This is exactly the same set as before (why?), so we can rewrite (3.1) as

$$\{n + 3 : n \in \mathbf{N}, 0 \leq n \leq 5\} = \{8 - m : m \in \mathbf{N}, 0 \leq m \leq 5\}.$$

Now it is easy to see (using (3.1.4)) why this identity is true: every number of the form $n + 3$, where n is a natural number between 0 and 5, is also of the form $8 - m$ where $m := 5 - n$ (note that m is therefore also a natural number between 0 and 5); conversely, every number of the form $8 - m$, where n is a natural number between 0 and 5, is also of the form $n + 3$, where $n := 5 - m$ (note that n is therefore a natural number between 0 and 5). Observe how much more confusing the above explanation of (3.1) would have been if we had not changed one of the n 's to an m first!

Exercise 3.1.1. Show that the definition of equality in (3.1.4) is reflexive, symmetric, and transitive.

Exercise 3.1.2. Using only Definition 3.1.4, Axiom 3.2, and Axiom 3.3, prove that the sets \emptyset , $\{\emptyset\}$, $\{\{\emptyset\}\}$, and $\{\emptyset, \{\emptyset\}\}$ are all distinct (i.e., no two of them are equal to each other).

Exercise 3.1.3. Prove the remaining claims in Lemma 3.1.12.

Exercise 3.1.4. Prove the remaining claims in Proposition 3.1.17. (Hint: one can use the first three claims to prove the fourth.)

Exercise 3.1.5. Let A, B be sets. Show that the three statements $A \subseteq B$, $A \cup B = B$, $A \cap B = A$ are logically equivalent (any one of them implies the other two).

Exercise 3.1.6. Prove the remaining claims in Proposition 3.1.27. (Hint: one can use some of these claims to prove others. Some of the claims have also appeared previously in Lemma 3.1.12).

Exercise 3.1.7. Let A, B, C be sets. Show that $A \cap B \subseteq A$ and $A \cap B \subseteq B$. Furthermore, show that $C \subseteq A$ and $C \subseteq B$ if and

only if $C \subseteq A \cap B$. In a similar spirit, show that $A \subseteq A \cup B$ and $B \subseteq A \cup B$, and furthermore that $A \subseteq C$ and $B \subseteq C$ if and only if $A \cup B \subseteq C$.

Exercise 3.1.8. Let A, B be sets. Prove the *absorption laws* $A \cap (A \cup B) = A$ and $A \cup (A \cap B) = A$.

Exercise 3.1.9. Let A, B, X be sets such that $A \cup B = X$ and $A \cap B = \emptyset$. Show that $A = X \setminus B$ and $B = X \setminus A$.

Exercise 3.1.10. Let A and B be sets. Show that the three sets $A \setminus B$, $A \cap B$, and $B \setminus A$ are disjoint, and that their union is $A \cup B$.

Exercise 3.1.11. Show that the axiom of replacement implies the axiom of specification.

3.2 Russell's paradox (Optional)

Many of the axioms introduced in the previous section have a similar flavor: they both allow us to form a set consisting of all the elements which have a certain property. They are both plausible, but one might think that they could be unified, for instance by introducing the following axiom:

Axiom 3.8 (Universal specification). (*Dangerous!*) Suppose for every object x we have a property $P(x)$ pertaining to x (so that for every x , $P(x)$ is either a true statement or a false statement). Then there exists a set $\{x : P(x) \text{ is true}\}$ such that for every object y ,

$$y \in \{x : P(x) \text{ is true}\} \iff P(y) \text{ is true.}$$

This axiom is also known as the *axiom of comprehension*. It asserts that every property corresponds to a set; if we assumed that axiom, we could talk about the set of all blue objects, the set of all natural numbers, the set of all sets, and so forth. This axiom also implies most of the axioms in the previous section (Exercise 3.2.1). Unfortunately, this axiom cannot be introduced into set theory, because it creates a logical contradiction known as *Russell's paradox*, discovered by the philosopher and logician Bertrand Russell

(1872–1970) in 1901. The paradox runs as follows. Let $P(x)$ be the statement

$$P(x) \iff \text{“}x \text{ is a set, and } x \notin x\text{”};$$

i.e., $P(x)$ is true only when x is a set which does not contain itself. For instance, $P(\{2, 3, 4\})$ is true, since the set $\{2, 3, 4\}$ is not one of the three elements 2, 3, 4 of $\{2, 3, 4\}$. On the other hand, if we let S be the set of all sets (which would we know to exist from the axiom of universal specification), then since S is itself a set, it is an element of S , and so $P(S)$ is true. Now use the axiom of universal specification to create the set

$$\Omega := \{x : P(x) \text{ is true}\} = \{x : x \text{ is a set and } x \notin x\},$$

i.e., the set of all sets which do not contain themselves. Now ask the question: does Ω contain itself, i.e. is $\Omega \in \Omega$? If Ω did contain itself, then by definition this means that $P(\Omega)$ is true, i.e., Ω is a set and $\Omega \notin \Omega$. On the other hand, if Ω did not contain itself, then $P(\Omega)$ would be true, and hence $\Omega \in \Omega$. Thus in either case we have both $\Omega \in \Omega$ and $\Omega \notin \Omega$, which is absurd.

The problem with the above axiom is that it creates sets which are far too “large” - for instance, we can use that axiom to talk about the set of *all* objects (a so-called “universal set”). Since sets are themselves objects (Axiom 3.1), this means that sets are allowed to contain themselves, which is a somewhat silly state of affairs. One way to informally resolve this issue is to think of objects as being arranged in a hierarchy. At the bottom of the hierarchy are the *primitive objects* - the objects that are not sets¹, such as the natural number 37. Then on the next rung of the hierarchy there are sets whose elements consist only of primitive objects, such as $\{3, 4, 7\}$ or the empty set \emptyset , let's call these “primitive sets” for now. Then there are sets whose elements consist only of primitive objects and primitive sets, such as $\{3, 4, 7, \{3, 4, 7\}\}$. Then we can form sets out of these objects, and so forth. The

¹In pure set theory, there will be no primitive objects, but there will be one primitive set \emptyset on the next rung of the hierarchy.

point is that at each stage of the hierarchy we only see sets whose elements consist of objects at lower stages of the hierarchy, and so at no stage do we ever construct a set which contains itself.

To actually formalize the above intuition of a hierarchy of objects is actually rather complicated, and we will not do so here. Instead, we shall simply postulate an axiom which ensures that absurdities such as Russell's paradox do not appear.

Axiom 3.9 (Regularity). *If A is a non-empty set, then there is at least one element x of A which is either not a set, or is disjoint from A .*

The point of this axiom (which is also known as the *axiom of foundation*) is that it is asserting that at least one of the elements of A is so low on the hierarchy of objects that it does not contain any of the other elements of A . For instance, if $A = \{\{3, 4\}, \{3, 4, \{3, 4\}\}\}$, then the element $\{3, 4\} \in A$ does not contain any of the elements of A (neither 3 nor 4 lies in A), although the element $\{3, 4, \{3, 4\}\}$, being somewhat higher in the hierarchy, does contain an element of A , namely $\{3, 4\}$. One particular consequence of this axiom is that sets are no longer allowed to contain themselves (Exercise 3.2.2).

One can legitimately ask whether we really need this axiom in our set theory, as it is certainly less intuitive than our other axioms. For the purposes of doing analysis, it turns out in fact that this axiom is never needed; all the sets we consider in analysis are typically very low on the hierarchy of objects, for instance being sets of primitive objects, or sets of sets of primitive objects, or at worst sets of sets of sets of primitive objects. However it is necessary to include this axiom in order to perform more advanced set theory, and so we have included this axiom in the text (but in an optional section) for sake of completeness.

Exercise 3.2.1. Show that the universal specification axiom, Axiom 3.8, if assumed to be true, would imply Axioms 3.2, 3.3, 3.4, 3.5, and 3.6. (If we assume that all natural numbers are objects, we also obtain Axiom 3.7). Thus, this axiom, if permitted, would simplify the foundations of set theory tremendously (and can be

viewed as one basis for an intuitive model of set theory known as “naive set theory”). Unfortunately, as we have seen, Axiom 3.8 is “too good to be true”!

Exercise 3.2.2. Use the axiom of regularity (and the singleton set axiom) to show that if A is a set, then $A \notin A$. Furthermore, show that if A and B are two sets, then either $A \notin B$ or $B \notin A$ (or both).

Exercise 3.2.3. Show (assuming the other axioms of set theory) that the universal specification axiom, Axiom 3.8, is equivalent to an axiom postulating the existence of a “universal set” Ω consisting of all objects (i.e., for all objects x , we have $x \in \Omega$). In other words, if Axiom 3.8 is true, then a universal set exists, and conversely, if a universal set exists, then Axiom 3.8 is true. (This may explain why Axiom 3.8 is called the axiom of *universal* specification). Note that if a universal set Ω existed, then we would have $\Omega \in \Omega$ by Axiom 3.1, contradicting Exercise 3.2.2. Thus the axiom of foundation specifically rules out the axiom of universal specification.

3.3 Functions

In order to do analysis, it is not particularly useful to just have the notion of a set; we also need the notion of a *function* from one set to another. Informally, a function $f : X \rightarrow Y$ from one set X to another set Y is an operation which assigns each element (or “input”) x in X , a single element (or “output”) $f(x)$ in Y ; we have already used this informal concept in the previous chapter when we discussed the natural numbers. The formal definition is as follows.

Definition 3.3.1 (Functions). Let X, Y be sets, and let $P(x, y)$ be a property pertaining to an object $x \in X$ and an object $y \in Y$, such that for every $x \in X$, there is exactly one $y \in Y$ for which $P(x, y)$ is true (this is sometimes known as the *vertical line test*). Then we define the *function* $f : X \rightarrow Y$ defined by P on the domain X and range Y to be the object which, given any input

$x \in X$, assigns an output $f(x) \in Y$, defined to be the unique object $f(x)$ for which $P(x, f(x))$ is true. Thus, for any $x \in X$ and $y \in Y$,

$$y = f(x) \iff P(x, y) \text{ is true.}$$

Functions are also referred to as *maps* or *transformations*, depending on the context. (They are also sometimes called *morphisms*, although to be more precise, a morphism refers to a more general class of object, which may or may not correspond to actual functions, depending on the context).

Example 3.3.2. Let $X = \mathbf{N}$, $Y = \mathbf{N}$, and let $P(x, y)$ be the property that $y = x++$. Then for each $x \in \mathbf{N}$ there is exactly one y for which $P(x, y)$ is true, namely $y = x++$. Thus we can define a function $f : \mathbf{N} \rightarrow \mathbf{N}$ associated to this property, so that $f(x) = x++$ for all x ; this is the *increment* function on \mathbf{N} , which takes a natural number as input and returns its increment as output. Thus for instance $f(4) = 5$, $f(2n + 3) = 2n + 4$ and so forth. One might also hope to define a *decrement* function $g : \mathbf{N} \rightarrow \mathbf{N}$ associated to the property $P(x, y)$ defined by $y++ = x$, i.e., $g(x)$ would be the number whose increment is x . Unfortunately this does not define a function, because when $x = 0$ there is no natural number y whose increment is equal to x (Axiom 2.3). On the other hand, we can legitimately define a decrement function $h : \mathbf{N} \setminus \{0\} \rightarrow \mathbf{N}$ associated to the property $P(x, y)$ defined by $y++ = x$, because when $x \in \mathbf{N} \setminus \{0\}$ there is indeed exactly one natural number y such that $y++ = x$, thanks to Lemma 2.2.10. Thus for instance $h(4) = 3$ and $h(2n + 3) = h(2n + 2)$, but $h(0)$ is undefined since 0 is not in the domain $\mathbf{N} \setminus \{0\}$.

Example 3.3.3. (Informal) This example requires the real numbers \mathbf{R} , which we will define in Chapter 5. One could try to define a square root function $\sqrt{\cdot} : \mathbf{R} \rightarrow \mathbf{R}$ by associating it to the property $P(x, y)$ defined by $y^2 = x$, i.e., we would want \sqrt{x} to be the number y such that $y^2 = x$. Unfortunately there are two problems which prohibit this definition from actually creating a function. The first is that there exist real numbers x for which $P(x, y)$ is

never true, for instance if $x = -1$ then there is no real number y such that $y^2 = x$. This problem however can be solved by restricting the domain from \mathbf{R} to the right half-line $[0, +\infty)$. The second problem is that even when $x \in [0, +\infty)$, it is possible for there to be more than one y in the range \mathbf{R} for which $y^2 = x$, for instance if $x = 4$ then both $y = 2$ and $y = -2$ obey the property $P(x, y)$, i.e., both $+2$ and -2 are square roots of 4. This problem can however be solved by restricting the range of \mathbf{R} to $[0, +\infty)$. Once one does this, then one can correctly define a square root function $\sqrt{\cdot} : [0, +\infty) \rightarrow [0, +\infty)$ using the relation $y^2 = x$, thus \sqrt{x} is the unique number $y \in [0, +\infty)$ such that $y^2 = x$.

One common way to define a function is simply to specify its domain, its range, and how one generates the output $f(x)$ from each input; this is known as an *explicit* definition of a function. For instance, the function f in Example 3.3.2 could be defined explicitly by saying that f has domain and range equal to \mathbf{N} , and $f(x) := x++$ for all $x \in \mathbf{N}$. In other cases we only define a function f by specifying what property $P(x, y)$ links the input x with the output $f(x)$; this is an *implicit* definition of a function. For instance, the square root function \sqrt{x} in Example 3.3.3 was defined implicitly by the relation $(\sqrt{x})^2 = x$. Note that an implicit definition is only valid if we know that for every input there is exactly one output which obeys the implicit relation. In many cases we omit specifying the domain and range of a function for brevity, and thus for instance we could refer to the function f in Example 3.3.2 as “the function $f(x) := x++$ ”, “the function $x \mapsto x++$ ”, “the function $x++$ ”, or even the extremely abbreviated “++”. However, too much of this abbreviation can be dangerous; sometimes it is important to know what the domain and range of the function is.

We observe that functions obey the axiom of substitution: if $x = x'$, then $f(x) = f(x')$ (why?). In other words, equal inputs imply equal outputs. On the other hand, unequal inputs do not necessarily ensure unequal outputs, as the following example shows:

Example 3.3.4. Let $X = \mathbf{N}$, $Y = \mathbf{N}$, and let $P(x, y)$ be the property that $y = 7$. Then certainly for every $x \in \mathbf{N}$ there is exactly one y for which $P(x, y)$ is true, namely the number 7. Thus we can create a function $f : \mathbf{N} \rightarrow \mathbf{N}$ associated to this property; it is simply the *constant function* which assigns the output of $f(x) = 7$ to each input $x \in \mathbf{N}$. Thus it is certainly possible for different inputs to generate the same output.

Remark 3.3.5. We are now using parentheses $()$ to denote several different things in mathematics; on one hand, we are using them to clarify the order of operations (compare for instance $2 + (3 \times 4) = 14$ with $(2 + 3) \times 4 = 20$), but on the other hand we also use parentheses to enclose the argument $f(x)$ of a function or of a property such as $P(x)$. However, the two usages of parentheses usually are unambiguous from context. For instance, if a is a number, then $a(b + c)$ denotes the expression $a \times (b + c)$, whereas if f is a function, then $f(b + c)$ denotes the output of f when the input is $b + c$. Sometimes the argument of a function is denoted by subscripting instead of parentheses; for instance, a sequence of natural numbers $a_0, a_1, a_2, a_3, \dots$ is, strictly speaking, a function from \mathbf{N} to \mathbf{N} , but is denoted by $n \mapsto a_n$ rather than $n \mapsto a(n)$.

Remark 3.3.6. Strictly speaking, functions are not sets, and sets are not functions; it does not make sense to ask whether an object x is an element of a function f , and it does not make sense to apply a set A to an input x to create an output $A(x)$. On the other hand, it is possible to start with a function $f : X \rightarrow Y$ and construct its *graph* $\{(x, f(x)) : x \in X\}$, which describes the function completely: see Section 3.5.

We now define some basic concepts and notions for functions. The first notion is that of equality.

Definition 3.3.7 (Equality of functions). Two functions $f : X \rightarrow Y$, $g : X \rightarrow Y$ with the same domain and range are said to be *equal*, $f = g$, if and only if $f(x) = g(x)$ for *all* $x \in X$. (If $f(x)$ and $g(x)$ agree for some values of x , but not others, then we do

not consider f and g to be equal².)

Example 3.3.8. The functions $x \mapsto x^2 + 2x + 1$ and $x \mapsto (x + 1)^2$ are equal on the domain \mathbf{R} . The functions $x \mapsto x$ and $x \mapsto |x|$ are equal on the positive real axis, but are not equal on \mathbf{R} ; thus the concept of equality of functions can depend on the choice of domain.

Example 3.3.9. A rather boring example of a function is the *empty function* $f : \emptyset \rightarrow X$ from the empty set to an arbitrary set X . Since the empty set has no elements, we do not need to specify what f does to any input. Nevertheless, just as the empty set is a set, the empty function is a function, albeit not a particularly interesting one. Note that for each set X , there is only one function from \emptyset to X , since Definition 3.3.7 asserts that all functions from \emptyset to X are equal (why?).

This notion of equality obeys the usual axioms (Exercise 3.3.1).

A fundamental operation available for functions is *composition*.

Definition 3.3.10 (Composition). Let $f : X \rightarrow Y$ and $g : Y \rightarrow Z$ be two functions, such that the range of f is the same set as the domain of g . We then define the *composition* $g \circ f : X \rightarrow Z$ of the two functions g and f to be the function defined explicitly by the formula

$$(g \circ f)(x) := g(f(x)).$$

If the range of f does not match the domain of g , we leave the composition $g \circ f$ undefined.

It is easy to check that composition obeys the axiom of substitution (Exercise 3.3.1).

Example 3.3.11. Let $f : \mathbf{N} \rightarrow \mathbf{N}$ be the function $f(n) := 2n$, and let $g : \mathbf{N} \rightarrow \mathbf{N}$ be the function $g(n) := n + 3$. Then $g \circ f$ is the function

$$g \circ f(n) = g(f(n)) = g(2n) = 2n + 3,$$

²Later on in this text, we shall introduce a weaker notion of equality, that of two functions being *equal almost everywhere*. However, we will not encounter this slightly different notion for a while yet.

thus for instance $g \circ f(1) = 5$, $g \circ f(2) = 7$, and so forth. Meanwhile, $f \circ g$ is the function

$$f \circ g(n) = f(g(n)) = f(n + 3) = 2(n + 3) = 2n + 6,$$

thus for instance $f \circ g(1) = 8$, $f \circ g(2) = 10$, and so forth.

The above example shows that composition is not commutative: $f \circ g$ and $g \circ f$ are not necessarily the same function. However, composition is still associative:

Lemma 3.3.12 (Composition is associative). *Let $f : X \rightarrow Y$, $g : Y \rightarrow Z$, and $h : Z \rightarrow W$ be functions. Then $f \circ (g \circ h) = (f \circ g) \circ h$.*

Proof. Since $g \circ h$ is a function from Y to W , $f \circ (g \circ h)$ is a function from X to W . Similarly $f \circ g$ is a function from X to Z , and hence $(f \circ g) \circ h$ is a function from X to W . Thus $f \circ (g \circ h)$ and $(f \circ g) \circ h$ have the same domain and range. In order to check that they are equal, we see from Definition 3.3.7 that we have to verify that $(f \circ (g \circ h))(x) = ((f \circ g) \circ h)(x)$ for all $x \in X$. But by Definition 3.3.10

$$\begin{aligned} (f \circ (g \circ h))(x) &= f((g \circ h)(x)) \\ &= f(g(h(x))) \\ &= (f \circ g)(h(x)) \\ &= ((f \circ g) \circ h)(x) \end{aligned}$$

as desired. □

Remark 3.3.13. Note that while g appears to the left of f in the expression $g \circ f$, the function $g \circ f$ applies the right-most function f first, before applying g . This is often confusing at first; it arises because we traditionally place a function f to the left of its input x rather than to the right. (There are some alternate mathematical notations in which the function is placed to the right of the input, thus we would write xf instead of $f(x)$, but this notation has often proven to be more confusing than clarifying, and has not as yet become particularly popular.)

We now describe certain special types of functions: *one-to-one* functions, *onto* functions, and *invertible* functions.

Definition 3.3.14 (One-to-one functions). A function f is *one-to-one* (or *injective*) if different elements map to different elements:

$$x \neq x' \implies f(x) \neq f(x').$$

Equivalently, a function is one-to-one if

$$f(x) = f(x') \implies x = x'.$$

Example 3.3.15. (Informal) The function $f : \mathbf{Z} \rightarrow \mathbf{Z}$ defined by $f(n) := n^2$ is not one-to-one because the distinct elements -1 , 1 map to the same element 1 . On the other hand, if we restrict this function to the natural numbers, defining the function $g : \mathbf{N} \rightarrow \mathbf{Z}$ by $g(n) := n^2$, then g is now a one-to-one function. Thus the notion of a one-to-one function depends not just on what the function does, but also what its domain is.

Remark 3.3.16. If a function $f : X \rightarrow Y$ is not one-to-one, then one can find distinct x and x' in the domain X such that $f(x) = f(x')$, thus one can find two inputs which map to one output. Because of this, we say that f is *two-to-one* instead of *one-to-one*.

Definition 3.3.17 (Onto functions). A function f is *onto* (or *surjective*) if $f(X) = Y$, i.e., every element in Y comes from applying f to some element in X :

For every $y \in Y$, there exists $x \in X$ such that $f(x) = y$.

Example 3.3.18. (Informal) The function $f : \mathbf{Z} \rightarrow \mathbf{Z}$ defined by $f(n) := n^2$ is not onto because the negative numbers are not in the image of f . However, if we restrict the range \mathbf{Z} to the set $A := \{n^2 : n \in \mathbf{Z}\}$ of square numbers, then the function $g : \mathbf{Z} \rightarrow A$ defined by $g(n) := n^2$ is now onto. Thus the notion of an onto function depends not just on what the function does, but also what its range is.

Remark 3.3.19. The concepts of injectivity and surjectivity are in many ways dual to each other; see Exercises 3.3.2, 3.3.4, 3.3.5 for some evidence of this.

Definition 3.3.20 (Bijective functions). Functions $f : X \rightarrow Y$ which are both one-to-one and onto are also called *bijective* or *invertible*.

Example 3.3.21. Let $f : \{0, 1, 2\} \rightarrow \{3, 4\}$ be the function $f(0) := 3$, $f(1) := 3$, $f(2) := 4$. This function is not bijective because if we set $y = 3$, then there is more than one x in $\{0, 1, 2\}$ such that $f(x) = y$ (this is a failure of injectivity). Now let $g : \{0, 1\} \rightarrow \{2, 3, 4\}$ be the function $g(0) := 2$, $g(1) := 3$; then g is not bijective because if we set $y = 4$, then there is no x for which $g(x) = y$ (this is a failure of surjectivity). Now let $h : \{0, 1, 2\} \rightarrow \{3, 4, 5\}$ be the function $h(0) := 3$, $h(1) := 4$, $h(2) := 5$. Then h is bijective, because each of the elements 3, 4, 5 comes from exactly one element from 0, 1, 2.

Example 3.3.22. The function $f : \mathbf{N} \rightarrow \mathbf{N} \setminus \{0\}$ defined by $f(n) := n++$ is a bijection (in fact, this fact is simply restating Axioms 2.2, 2.3, 2.4). On the other hand, the function $g : \mathbf{N} \rightarrow \mathbf{N}$ defined by the same definition $g(n) := n++$ is not a bijection. Thus the notion of a bijective function depends not just on what the function does, but also what its range (and domain) are.

Remark 3.3.23. If a function $x \mapsto f(x)$ is bijective, then we sometimes call f a *perfect matching* or *one-to-one correspondence* (not to be confused with the notion of a one-to-one function), and denote the action of f using the notation $x \leftrightarrow f(x)$ instead of $x \mapsto f(x)$. Thus for instance the function h in the above example is the one-to-one correspondence $0 \leftrightarrow 3$, $1 \leftrightarrow 4$, $2 \leftrightarrow 5$.

Remark 3.3.24. A common error is to say that a function $f : X \rightarrow Y$ is bijective iff “for every x in X , there is exactly one y in Y such that $y = f(x)$.” This is not what it means for f to be bijective; it is what it means for f to be a *function*: each input gives exactly one output. A function cannot map one element to

two different elements, for instance one cannot have a function f for which $f(0) = 1$ and also $f(0) = 2$. The functions f, g given in the previous example are not bijective, but they are still functions, since each input still gives exactly one output.

If f is bijective, then for every $y \in Y$, there is exactly one x such that $f(x) = y$ (there is at least one because of surjectivity, and at most one because of injectivity). This value of x is denoted $f^{-1}(y)$; thus f^{-1} is a function from Y to X .

Exercise 3.3.1. Show that the definition of equality in Definition 3.3.7 is reflexive, symmetric, and transitive. Also verify the substitution property: if $f, \tilde{f} : X \rightarrow Y$ and $g, \tilde{g} : Y \rightarrow Z$ are functions such that $f = \tilde{f}$ and $g = \tilde{g}$, then $f \circ g = \tilde{f} \circ \tilde{g}$.

Exercise 3.3.2. Let $f : X \rightarrow Y$ and $g : Y \rightarrow Z$ be functions. Show that if f and g are both injective, then so is $g \circ f$; similarly, show that if f and g are both surjective, then so is $g \circ f$.

Exercise 3.3.3. When is the empty function injective? surjective? bijective?

Exercise 3.3.4. In this section we give some cancellation laws for composition. Let $f : X \rightarrow Y$, $\tilde{f} : X \rightarrow Y$, $g : Y \rightarrow Z$, and $\tilde{g} : Y \rightarrow Z$ be functions. Show that if $g \circ f = g \circ \tilde{f}$ and g is injective, then $f = \tilde{f}$. Is the same statement true if g is not injective? Show that if $g \circ f = \tilde{g} \circ f$ and f is surjective, then $g = \tilde{g}$. Is the same statement true if f is not surjective?

Exercise 3.3.5. Let $f : X \rightarrow Y$ and $g : Y \rightarrow Z$ be functions. Show that if $g \circ f$ is injective, then f must be injective. Is it true that g must also be injective? Show that if $g \circ f$ is surjective, then g must be surjective. Is it true that f must also be surjective?

Exercise 3.3.6. Let $f : X \rightarrow Y$ be a bijective function, and let $f^{-1} : Y \rightarrow X$ be its inverse. Verify the cancellation laws $f^{-1}(f(x)) = x$ for all $x \in X$ and $f(f^{-1}(y)) = y$ for all $y \in Y$. Conclude that f^{-1} is also invertible, and has f as its inverse (thus $(f^{-1})^{-1} = f$).

Exercise 3.3.7. Let $f : X \rightarrow Y$ and $g : Y \rightarrow Z$ be functions. Show that if f and g are bijective, then so is $g \circ f$, and we have $(g \circ f)^{-1} = f^{-1} \circ g^{-1}$.

Exercise 3.3.8. If X is a subset of Y , let $\iota_{X \rightarrow Y} : X \rightarrow Y$ be the *inclusion map from X to Y* , defined by mapping $x \mapsto x$ for all $x \in X$, i.e., $\iota_{X \rightarrow Y}(x) := x$ for all $x \in X$. The map $\iota_{X \rightarrow X}$ is in particular called the *identity map on X* .

- Show that if $X \subseteq Y \subseteq Z$ then $\iota_{Y \rightarrow Z} \circ \iota_{X \rightarrow Y} = \iota_{X \rightarrow Z}$.
- Show that if $f : A \rightarrow B$ is any function, then $f = f \circ \iota_{A \rightarrow A} = \iota_{B \rightarrow B} \circ f$.
- Show that, if $f : A \rightarrow B$ is a bijective function, then $f \circ f^{-1} = \iota_{B \rightarrow B}$ and $f^{-1} \circ f = \iota_{A \rightarrow A}$.
- Show that if X and Y are disjoint sets, and $f : X \rightarrow Z$ and $g : Y \rightarrow Z$ are functions, then there is a unique function $h : X \cup Y \rightarrow Z$ such that $h \circ \iota_{X \rightarrow X \cup Y} = f$ and $h \circ \iota_{Y \rightarrow X \cup Y} = g$.

3.4 Images and inverse images

We know that a function $f : X \rightarrow Y$ from a set X to a set Y can take individual elements $x \in X$ to elements $f(x) \in Y$. Functions can also take subsets in X to subsets in Y :

Definition 3.4.1 (Images of sets). If $f : X \rightarrow Y$ is a function from X to Y , and S is a set in X , we define $f(S)$ to be the set

$$f(S) := \{f(x) : x \in S\};$$

this set is a subset of Y , and is sometimes called the *image* of S under the map f . We sometimes call $f(S)$ the *forward image* of S to distinguish it from the concept of the *inverse image* $f^{-1}(S)$ of S , which is defined below.

Note that the set $f(S)$ is well-defined thanks to the axiom of replacement (Axiom 3.6). One can also define $f(S)$ using the axiom of specification (Axiom 3.5) instead of replacement, but we leave this as a challenge to the reader.

Example 3.4.2. If $f : \mathbf{N} \rightarrow \mathbf{N}$ is the map $f(x) = 2x$, then the forward image of $\{1, 2, 3\}$ is $\{2, 4, 6\}$:

$$f(\{1, 2, 3\}) = \{2, 4, 6\}.$$

More informally, to compute $f(S)$, we take every element x of S , and apply f to each element individually, and then put all the resulting objects together to form a new set.

In the above example, the image had the same size as the original set. But sometimes the image can be smaller, because f is not one-to-one (see Definition 3.3.14):

Example 3.4.3. (Informal) Let \mathbf{Z} be the set of integers (which we will define rigorously in the next section) and let $f : \mathbf{Z} \rightarrow \mathbf{Z}$ be the map $f(x) = x^2$, then

$$f(\{-1, 0, 1, 2\}) = \{0, 1, 4\}.$$

Note that f is not one-to-one because $f(-1) = f(1)$.

Note that

$$x \in S \implies f(x) \in f(S)$$

but in general

$$f(x) \in f(S) \not\Rightarrow x \in S;$$

for instance in the above informal example, $f(-2)$ is in $f(\{-1, 0, 1, 2\})$, but -2 is not in $\{-1, 0, 1, 2\}$. The correct statement is

$$y \in f(S) \iff y = f(x) \text{ for some } x \in S$$

(why?).

Definition 3.4.4 (Inverse images). If U is a subset of Y , we define the set $f^{-1}(U)$ to be the set

$$f^{-1}(U) := \{x \in X : f(x) \in U\}.$$

In other words, $f^{-1}(U)$ consists of all the stuff in X which maps into U :

$$f(x) \in U \iff x \in f^{-1}(U).$$

We call $f^{-1}(U)$ the *inverse image* of U .

Example 3.4.5. If $f : \mathbf{N} \rightarrow \mathbf{N}$ is the map $f(x) = 2x$, then $f(\{1, 2, 3\}) = \{2, 4, 6\}$, but $f^{-1}(\{1, 2, 3\}) = \{1\}$. Thus the forward image of $\{1, 2, 3\}$ and the backwards image of $\{1, 2, 3\}$ are quite different sets. Also note that

$$f(f^{-1}(\{1, 2, 3\})) \neq \{1, 2, 3\}$$

(why?).

Example 3.4.6. (Informal) If $f : \mathbf{Z} \rightarrow \mathbf{Z}$ is the map $f(x) = x^2$, then

$$f^{-1}(\{0, 1, 4\}) = \{-2, -1, 0, 1, 2\}.$$

Note that f does not have to be invertible in order for $f^{-1}(U)$ to make sense. Also note that images and inverse images do not quite invert each other, for instance we have

$$f^{-1}(f(\{-1, 0, 1, 2\})) \neq \{-1, 0, 1, 2\}$$

(why?).

Note that we have now defined f^{-1} in two slightly different ways, but this is not an issue because both definitions are equivalent (Exercise 3.4.1).

As remarked earlier, functions are not sets. However, we do consider functions to be a type of object, and in particular we should be able to consider sets of functions. In particular, we should be able to consider the set of *all* functions from a set X to a set Y . To do this we need to introduce another axiom to set theory:

Axiom 3.10 (Power set axiom). *Let X and Y be sets. Then there exists a set, denoted Y^X , which consists of all the functions from X to Y , thus*

$$f \in Y^X \iff (f \text{ is a function with domain } X \text{ and range } Y).$$

Example 3.4.7. Let $X = \{4, 7\}$ and $Y = \{0, 1\}$. Then the set Y^X consists of four functions: the function that maps $4 \mapsto 0$ and

$7 \mapsto 0$; the function that maps $4 \mapsto 0$ and $7 \mapsto 1$; the function that maps $4 \mapsto 1$ and $7 \mapsto 0$; and the function that maps $4 \mapsto 1$ and $7 \mapsto 1$. The reason we use the notation Y^X to denote this set is that if Y has n elements and X has m elements, then one can show that Y^X has n^m elements; see Proposition 3.6.13(f).

One consequence of this axiom is

Lemma 3.4.8. *Let X be a set. Then the set*

$$\{Y : Y \text{ is a subset of } X\}$$

is a set.

Proof. See Exercise 3.4.6. □

Remark 3.4.9. The set $\{Y : Y \text{ is a subset of } X\}$ is known as the *power set* of X and is denoted 2^X . For instance, if a, b, c are distinct objects, we have

$$2^{\{a,b,c\}} = \{\emptyset, \{a\}, \{b\}, \{c\}, \{a, b\}, \{a, c\}, \{b, c\}, \{a, b, c\}\}.$$

Note that while $\{a, b, c\}$ has 3 elements, $2^{\{a,b,c\}}$ has $2^3 = 8$ elements. This gives a hint as to why we refer to the power set of X as 2^X ; we return to this issue in Chapter 8.

For sake of completeness, let us now add one further axiom to our set theory, in which we enhance the axiom of pairwise union to allow unions of much larger collections of sets.

Axiom 3.11 (Union). *Let A be a set, all of whose elements are themselves sets. Then there exists a set $\bigcup A$ whose elements are precisely those objects which are elements of the elements of A , thus for all objects x*

$$x \in \bigcup A \iff (x \in S \text{ for some } S \in A).$$

Example 3.4.10. If $A = \{\{2, 3\}, \{3, 4\}, \{4, 5\}\}$, then $\bigcup A = \{2, 3, 4, 5\}$ (why?).

The axiom of union, combined with the axiom of pair set, implies the axiom of pairwise union (Exercise 3.4.8). Another important consequence of this axiom is that if one has some set I , and for every element $\alpha \in I$ we have some set A_α , then we can form the union set $\bigcup_{\alpha \in I} A_\alpha$ by defining

$$\bigcup_{\alpha \in I} A_\alpha := \bigcup \{A_\alpha : \alpha \in I\},$$

which is a set thanks to the axiom of replacement and the axiom of union. Thus for instance, if $I = \{1, 2, 3\}$, $A_1 := \{2, 3\}$, $A_2 := \{3, 4\}$, and $A_3 := \{4, 5\}$, then $\bigcup_{\alpha \in \{1, 2, 3\}} A_\alpha = \{2, 3, 4, 5\}$. More generally, we see that for any object y ,

$$y \in \bigcup_{\alpha \in I} A_\alpha \iff (y \in A_\alpha \text{ for some } \alpha \in I). \quad (3.2)$$

In situations like this, we often refer to I as an *index set*, and the elements α of this index set as *labels*; the sets A_α are then called a *family of sets*, and are *indexed* by the labels $\alpha \in I$. Note that if I was empty, then $\bigcup_{\alpha \in I} A_\alpha$ would automatically also be empty (why?).

We can similarly form intersections of families of sets, as long as the index set is non-empty. More specifically, given any non-empty set I , and given an assignment of a set A_α to each $\alpha \in I$, we can define the intersection $\bigcap_{\alpha \in I} A_\alpha$ by first choosing some element β of I (which we can do since I is non-empty), and setting

$$\bigcap_{\alpha \in I} A_\alpha := \{x \in A_\beta : x \in A_\alpha \text{ for all } \alpha \in I\}, \quad (3.3)$$

which is a set by the axiom of specification. This definition may look like it depends on the choice of β , but it does not (Exercise 3.4.9). Observe that for any object y ,

$$y \in \bigcap_{\alpha \in I} A_\alpha \iff (y \in A_\alpha \text{ for all } \alpha \in I) \quad (3.4)$$

(compare with (3.2)).

Remark 3.4.11. The axioms of set theory that we have introduced (Axioms 3.1-3.11, excluding the dangerous axiom Axiom 3.8) are known as the³ *Zermelo-Fraenkel axioms of set theory*. There is one further axiom we will eventually need, the famous *axiom of choice* (see Section 8.4), giving rise to the *Zermelo-Fraenkel-Choice (ZFC) axioms of set theory*, but we will not need this axiom for some time.

Exercise 3.4.1. Let $f : X \rightarrow Y$ be a bijective function, and let $f^{-1} : Y \rightarrow X$ be its inverse. Let V be any subset of Y . Prove that the forward image of V under f^{-1} is the same set as the inverse image of V under f ; thus the fact that both sets are denoted by $f^{-1}(V)$ will not lead to any inconsistency.

Exercise 3.4.2. Let $f : X \rightarrow Y$ be a function from one set X to another set Y , let S be a subset of X , and let U be a subset of Y . What, in general, can one say about $f^{-1}(f(S))$ and S ? What about $f(f^{-1}(U))$ and U ?

Exercise 3.4.3. Let A, B be two subsets of a set X , and let $f : X \rightarrow Y$ be a function. Show that $f(A \cap B) \subseteq f(A) \cap f(B)$, that $f(A) \setminus f(B) \subseteq f(A \setminus B)$, $f(A \cup B) = f(A) \cup f(B)$. For the first two statements, is it true that the \subseteq relation can be improved to $=$?

Exercise 3.4.4. Let $f : X \rightarrow Y$ be a function from one set X to another set Y , and let U, V be subsets of Y . Show that $f^{-1}(U \cup V) = f^{-1}(U) \cup f^{-1}(V)$, that $f^{-1}(U \cap V) = f^{-1}(U) \cap f^{-1}(V)$, and that $f^{-1}(U \setminus V) = f^{-1}(U) \setminus f^{-1}(V)$.

Exercise 3.4.5. Let $f : X \rightarrow Y$ be a function from one set X to another set Y . Show that $f(f^{-1}(S)) = S$ for every $S \subseteq Y$ if and only if f is surjective. Show that $f^{-1}(f(S)) = S$ for every $S \subseteq X$ if and only if f is injective.

Exercise 3.4.6. Prove Lemma 3.4.8. (Hint: Start with the set $\{0, 1\}^X$ and apply the replacement axiom, replacing each function f with $f^{-1}(\{1\})$.) This statement has a converse; see Exercise 3.5.11.

³These axioms are formulated slightly differently in other texts, but all the formulations can be shown to be equivalent to each other.

Exercise 3.4.7. Let X, Y be sets. Define a *partial function* from X to Y to be any function $f : X' \rightarrow Y'$ whose domain X' is a subset of X , and whose range Y' is a subset of Y . Show that the collection of all partial functions from X to Y is itself a set. (Hint: use Exercise 3.4.6, the power set axiom, the replacement axiom, and the union axiom.)

Exercise 3.4.8. Show that Axiom 3.4 can be deduced from Axiom 3.3 and Axiom 3.11.

Exercise 3.4.9. Show that if β and β' are two elements of a set I , and to each $\alpha \in I$ we assign a set A_α , then

$$\{x \in A_\beta : x \in A_\alpha \text{ for all } \alpha \in I\} = \{x \in A_{\beta'} : x \in A_\alpha \text{ for all } \alpha \in I\},$$

and so the definition of $\bigcap_{\alpha \in I} A_\alpha$ defined in (3.3) does not depend on β . Also explain why (3.4) is true.

Exercise 3.4.10. Suppose that I and J are two sets, and for all $\alpha \in I \cup J$ let A_α be a set. Show that $(\bigcup_{\alpha \in I} A_\alpha) \cup (\bigcup_{\alpha \in J} A_\alpha) = \bigcup_{\alpha \in I \cup J} A_\alpha$. If I and J are non-empty, show that $(\bigcap_{\alpha \in I} A_\alpha) \cap (\bigcap_{\alpha \in J} A_\alpha) = \bigcap_{\alpha \in I \cup J} A_\alpha$.

Exercise 3.4.11. Let X be a set, let I be a non-empty set, and for all $\alpha \in I$ let A_α be a subset of X . Show that

$$X \setminus \bigcup_{\alpha \in I} A_\alpha = \bigcap_{\alpha \in I} (X \setminus A_\alpha)$$

and

$$X \setminus \bigcap_{\alpha \in I} A_\alpha = \bigcup_{\alpha \in I} (X \setminus A_\alpha).$$

This should be compared with de Morgan's laws in Proposition 3.1.27 (although one cannot derive the above identities directly from de Morgan's laws, as I could be infinite).

3.5 Cartesian products

In addition to the basic operations of union, intersection, and differencing, another fundamental operation on sets is that of *Cartesian product*.

Definition 3.5.1 (Ordered pair). If x and y are any objects (possibly equal), we define the *ordered pair* (x, y) to be a new object, consisting of x as its first component and y as its second component. Two ordered pairs (x, y) and (x', y') are considered equal if and only if both their components match, i.e.

$$(x, y) = (x', y') \iff (x = x' \text{ and } y = y'). \quad (3.5)$$

This obeys the usual axioms of equality (Exercise 3.5.3). Thus for instance, the pair $(3, 5)$ is equal to the pair $(2 + 1, 3 + 2)$, but is distinct from the pairs $(5, 3)$, $(3, 3)$, and $(2, 5)$. (This is in contrast to sets, where $\{3, 5\}$ and $\{5, 3\}$ are equal).

Remark 3.5.2. Strictly speaking, this definition is partly an axiom, because we have simply postulated that given any two objects x and y , that an object of the form (x, y) exists. However, it is possible to define an ordered pair using the axioms of set theory in such a way that we do not need any further postulates (see Exercise 3.5.1).

Remark 3.5.3. We have now “overloaded” the parenthesis symbols $()$ once again; they now are not only used to denote grouping of operators and arguments of functions, but also to enclose ordered pairs. This is usually not a problem in practice as one can still determine what usage the symbols $()$ were intended for from context.

Definition 3.5.4 (Cartesian product). If X and Y are sets, then we define the *Cartesian product* $X \times Y$ to be the collection of ordered pairs, whose first component lies in X and second component lies in Y , thus

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

or equivalently

$$a \in (X \times Y) \iff (a = (x, y) \text{ for some } x \in X \text{ and } y \in Y).$$

Remark 3.5.5. We shall simply assume that our notion of ordered pair is such that whenever X and Y are sets, the Cartesian product $X \times Y$ is also a set. This is however not a problem in practice; see Exercise 3.5.1.

Example 3.5.6. If $X := \{1, 2\}$ and $Y := \{3, 4, 5\}$, then

$$X \times Y = \{(1, 3), (1, 4), (1, 5), (2, 3), (2, 4), (2, 5)\}$$

and

$$Y \times X = \{(3, 1), (4, 1), (5, 1), (3, 2), (4, 2), (5, 2)\}.$$

Thus, strictly speaking, $X \times Y$ and $Y \times X$ are different sets, although they are very similar. For instance, they always have the same number of elements (Exercise 3.6.5).

Let $f : X \times Y \rightarrow Z$ be a function whose domain $X \times Y$ is a Cartesian product of two other sets X and Y . Then f can either be thought of as a function of one variable, mapping the single input of an ordered pair (x, y) in $X \times Y$ to an output $f(x, y)$ in Z , or as a function of two variables, mapping an input $x \in X$ and another input $y \in Y$ to a single output $f(x, y)$ in Z . While the two notions are technically different, we will not bother to distinguish the two, and think of f simultaneously as a function of one variable with domain $X \times Y$ and as a function of two variables with domains X and Y . Thus for instance the addition operation $+$ on the natural numbers can now be re-interpreted as a function $+$: $\mathbf{N} \times \mathbf{N} \rightarrow \mathbf{N}$, defined by $(x, y) \mapsto x + y$.

One can of course generalize the concept of ordered pairs to ordered triples, ordered quadruples, etc:

Definition 3.5.7 (Ordered n -tuple and n -fold Cartesian product). Let n be a natural number. An *ordered n -tuple* $(x_i)_{1 \leq i \leq n}$ (also denoted (x_1, \dots, x_n)) is a collection of objects x_i , one for every natural number i between 1 and n ; we refer to x_i as the i^{th} *component* of the ordered n -tuple. Two ordered n -tuples $(x_i)_{1 \leq i \leq n}$ and $(y_i)_{1 \leq i \leq n}$ are said to be equal iff $x_i = y_i$ for all $1 \leq i \leq n$. If

$(X_i)_{1 \leq i \leq n}$ is an ordered n -tuple of sets, we define their *Cartesian product* $\prod_{1 \leq i \leq n} X_i$ (also denoted $\prod_{i=1}^n X_i$ or $X_1 \times \dots \times X_n$) by

$$\prod_{1 \leq i \leq n} X_i := \{(x_i)_{1 \leq i \leq n} : x_i \in X_i \text{ for all } 1 \leq i \leq n\}.$$

Again, this definition simply postulates that an ordered n -tuple and a Cartesian product always exist when needed, but using the axioms of set theory one can explicitly construct these objects (Exercise 3.5.2).

Remark 3.5.8. One can show that $\prod_{1 \leq i \leq n} X_i$ is indeed a set, by starting with the power set axiom to consider the set of all functions $i \mapsto x_i$ from the domain $\{1 \leq i \leq n\}$ to the range $\bigcup_{1 \leq i \leq n} X_i$, and then using the axiom of specification to restrict to those functions $i \mapsto x_i$ for which $x_i \in X_i$ for all $1 \leq i \leq n$. One can generalize this definition to infinite Cartesian products, see Definition 8.4.1.

Example 3.5.9. Let $a_1, b_1, a_2, b_2, a_3, b_3$ be objects, and let $X_1 := \{a_1, b_1\}$, $X_2 := \{a_2, b_2\}$, and $X_3 := \{a_3, b_3\}$. Then we have

$$\begin{aligned} X_1 \times X_2 \times X_3 &= \{(a_1, a_2, a_3), (a_1, a_2, b_3), (a_1, b_2, a_3), (a_1, b_2, b_3), \\ &\quad (b_1, a_2, a_3), (b_1, a_2, b_3), (b_1, b_2, a_3), (b_1, b_2, b_3)\} \\ (X_1 \times X_2) \times X_3 &= \\ &\quad \{((a_1, a_2), a_3), ((a_1, a_2), b_3), ((a_1, b_2), a_3), ((a_1, b_2), b_3), \\ &\quad ((b_1, a_2), a_3), ((b_1, a_2), b_3), ((b_1, b_2), a_3), ((b_1, b_2), b_3)\} \\ X_1 \times (X_2 \times X_3) &= \\ &\quad \{(a_1, (a_2, a_3)), (a_1, (a_2, b_3)), (a_1, (b_2, a_3)), (a_1, (b_2, b_3)), \\ &\quad (b_1, (a_2, a_3)), (b_1, (a_2, b_3)), (b_1, (b_2, a_3)), (b_1, (b_2, b_3))\}. \end{aligned}$$

Thus, strictly speaking, the sets $X_1 \times X_2 \times X_3$, $(X_1 \times X_2) \times X_3$, and $X_1 \times (X_2 \times X_3)$ are distinct. However, they are clearly very related to each other (for instance, there are obvious bijections between any two of the three sets), and it is common in practice to neglect the minor distinctions between these sets and pretend that they are in fact equal. Thus a function $f : X_1 \times X_2 \times X_3 \rightarrow Y$ can be

thought of as a function of one variable $(x_1, x_2, x_3) \in X_1 \times X_2 \times X_3$, or as a function of three variables $x_1 \in X_1, x_2 \in X_2, x_3 \in X_3$, or as a function of two variables $x_1 \in X_1, (x_2, x_3) \in X_3$, and so forth; we will not bother to distinguish between these different perspectives.

Remark 3.5.10. An ordered n -tuple x_1, \dots, x_n of objects is sometimes also called an *ordered sequence* of n elements, or a *finite sequence* for short. In Chapter 5 we shall also introduce the very useful concept of an *infinite sequence*.

Example 3.5.11. If x is an object, then (x) is a 1-tuple, which we shall identify with x itself (even though the two are, strictly speaking, not the same object). Then if X_1 is any set, then the Cartesian product $\prod_{1 \leq i \leq 1} X_i$ is just X_1 (why?). Also, the *empty Cartesian product* $\prod_{1 \leq i \leq 0} X_i$ gives, not the empty set $\{\}$, but rather the singleton set $\{()\}$ whose only element is the (*empty*) 0-tuple $()$.

If n is a natural number, we often write X^n as shorthand for the n -fold Cartesian product $X^n := \prod_{1 \leq i \leq n} X$. Thus X^1 is essentially the same set as X (if we ignore the distinction between an object x and the 1-tuple (x)), while X^2 is the Cartesian product $X \times X$. The set X^0 is a singleton set $\{()\}$ (why?).

We can now generalize the single choice lemma (Lemma 3.1.6) to allow for multiple (but finite) number of choices.

Lemma 3.5.12 (Finite choice). *Let $n \geq 1$ be a natural number, and for each natural number $1 \leq i \leq n$, let X_i be a non-empty set. Then there exists an n -tuple $(x_i)_{1 \leq i \leq n}$ such that $x_i \in X_i$ for all $1 \leq i \leq n$. In other words, if each X_i is non-empty, then the set $\prod_{1 \leq i \leq n} X_i$ is also non-empty.*

Proof. We induct on n (starting with the base case $n = 1$; the claim is also vacuously true with $n = 0$ but is not particularly interesting in that case). When $n = 1$ the claim follows from Lemma 3.1.6 (why?). Now suppose inductively that the claim has already been proven for some n ; we will now prove it for $n++$.

Let X_1, \dots, X_{n++} be a collection of non-empty sets. By induction hypothesis, we can find an n -tuple $(x_i)_{1 \leq i \leq n}$ such that $x_i \in X_i$ for all $1 \leq i \leq n$. Also, since X_{n++} is non-empty, by Lemma 3.1.6 we may find an object a such that $a \in X_{n++}$. If we thus define the $n++$ -tuple $(y_i)_{1 \leq i \leq n++}$ by setting $y_i := x_i$ when $1 \leq i \leq n$ and $y_i := a$ when $i = n++$ it is clear that $y_i \in X_i$ for all $1 \leq i \leq n++$, thus closing the induction. \square

Remark 3.5.13. It is intuitively plausible that this lemma should be extended to allow for an infinite number of choices, but this cannot be done automatically; it requires an additional axiom, the *axiom of choice*. See Section 8.4.

Exercise 3.5.1. Suppose we *define* the ordered pair (x, y) for any objects x and y by the formula $(x, y) := \{\{x\}, \{x, y\}\}$ (thus using several applications of Axiom 3.3). Thus for instance $(1, 2)$ is the set $\{\{1\}, \{1, 2\}\}$, $(2, 1)$ is the set $\{\{2\}, \{2, 1\}\}$, and $(1, 1)$ is the set $\{\{1\}\}$. Show that such a definition indeed obeys the property (3.5), and also whenever X and Y are sets, the Cartesian product $X \times Y$ is also a set. Thus this definition can be validly used as a definition of an ordered pair. For an additional challenge, show that the alternate definition $(x, y) := \{x, \{x, y\}\}$ also verifies (3.5) and is thus also an acceptable definition of ordered pair. (For this latter task one needs the axiom of regularity, and in particular Exercise 3.2.2.)

Exercise 3.5.2. Suppose we *define* an ordered n -tuple to be a surjective function $x : \{i \in \mathbf{N} : 1 \leq i \leq n\} \rightarrow X$ whose range is some arbitrary set X (so different ordered n -tuples are allowed to have different ranges); we then write x_i for $x(i)$, and also write x as $(x_i)_{1 \leq i \leq n}$. Using this definition, verify that we have $(x_i)_{1 \leq i \leq n} = (y_i)_{1 \leq i \leq n}$ if and only if $x_i = y_i$ for all $1 \leq i \leq n$. Also, show that if $(X_i)_{1 \leq i \leq n}$ are an ordered n -tuple of sets, then the Cartesian product, as defined in Definition 3.5.7, is indeed a set. (Hint: use Exercise 3.4.7 and the axiom of specification.)

Exercise 3.5.3. Show that the definitions of equality for ordered pair and ordered n -tuple obey the reflexivity, symmetry, and transitivity axioms.

Exercise 3.5.4. Let A, B, C be sets. Show that $A \times (B \cup C) = (A \times B) \cup (A \times C)$, that $A \times (B \cap C) = (A \times B) \cap (A \times C)$, and that $A \times (B \setminus C) = (A \times B) \setminus (A \times C)$. (One can of course prove similar identities in which the roles of the left and right factors of the Cartesian product are reversed.)

Exercise 3.5.5. Let A, B, C, D be sets. Show that $(A \times B) \cap (C \times D) = (A \cap C) \times (B \cap D)$. Is it true that $(A \times B) \cup (C \times D) = (A \cup C) \times (B \cup D)$? Is it true that $(A \times B) \setminus (C \times D) = (A \setminus C) \times (B \setminus D)$?

Exercise 3.5.6. Let A, B, C, D be non-empty sets. Show that $A \times B \subseteq C \times D$ if and only if $A \subseteq C$ and $B \subseteq D$, and that $A \times B = C \times D$ if and only if $A = C$ and $B = D$. What happens if the hypotheses that the A, B, C, D are all non-empty are removed?

Exercise 3.5.7. Let X, Y be sets, and let $\pi_{X \times Y \rightarrow X} : X \times Y \rightarrow X$ and $\pi_{X \times Y \rightarrow Y} : X \times Y \rightarrow Y$ be the maps defined by $\pi_{X \times Y \rightarrow X}(x, y) := x$ and $\pi_{X \times Y \rightarrow Y}(x, y) := y$. Show that for any functions $f : Z \rightarrow X$ and $g : Z \rightarrow Y$, there exists a unique function $h : Z \rightarrow X \times Y$ such that $\pi_{X \times Y \rightarrow X} \circ h = f$ and $\pi_{X \times Y \rightarrow Y} \circ h = g$. (Compare this to the last part of Exercise 3.3.8, and to Exercise 3.1.7.) This function h is known as the *direct sum* of f and g and is denoted $h = f \oplus g$.

Exercise 3.5.8. Let X_1, \dots, X_n be sets. Show that the Cartesian product $\prod_{i=1}^n X_i$ is empty if and only if at least one of the X_i is empty.

Exercise 3.5.9. Suppose that I and J are two sets, and for all $\alpha \in I$ let A_α be a set, and for all $\beta \in J$ let B_β be a set. Show that $(\bigcup_{\alpha \in I} A_\alpha) \cap (\bigcup_{\beta \in J} B_\beta) = \bigcup_{(\alpha, \beta) \in I \times J} (A_\alpha \cap B_\beta)$.

Exercise 3.5.10. If $f : X \rightarrow Y$ is a function, define the *graph* of f to be the subset of $X \times Y$ defined by $\{(x, f(x)) : x \in X\}$. Show that two functions $f : X \rightarrow Y$, $\tilde{f} : X \rightarrow Y$ are equal if and only if they have the same graph. Conversely, if G is any subset of $X \times Y$ with the property that for each $x \in X$, the set $\{y \in Y : (x, y) \in G\}$ has exactly one element (or in other words, G obeys the *vertical line test*), show that there is exactly one function $f : X \rightarrow Y$ whose graph is equal to G .

Exercise 3.5.11. Show that Axiom 3.10 can in fact be deduced from Lemma 3.4.8 and the other axioms of set theory, and thus

Lemma 3.4.8 can be used as an alternate formulation of the power set axiom. (Hint: for any two sets X and Y , use Lemma 3.4.8 and the axiom of specification to construct the set of all subsets of $X \times Y$ which obey the vertical line test. Then use Exercise 3.5.10 and the axiom of replacement.)

Exercise 3.5.12. The purpose of this exercise is to prove a rigorous version of Proposition 2.1.16. Let $f : \mathbf{N} \times \mathbf{N} \rightarrow \mathbf{N}$ be a function, and let c be a natural number. Show that there exists a function $a : \mathbf{N} \rightarrow \mathbf{N}$ such that

$$a(0) = c$$

and

$$a(n++) = f(n, a(n)) \text{ for all } n \in \mathbf{N},$$

and furthermore that this function is unique. (Hint: first show inductively, by a modification of the proof of Lemma 3.5.12, that for every natural number $N \in \mathbf{N}$, there exists a unique function $a_N : \{n \in \mathbf{N} : n \leq N\} \rightarrow \mathbf{N}$ such that $a_N(0) = c$ and $a_N(n++) = f(n, a(n))$ for all $n \in \mathbf{N}$ such that $n < N$.) For an additional challenge, prove this result without using any properties of the natural numbers other than the Peano axioms directly (in particular, without using the ordering of the natural numbers, and without appealing to Proposition 2.1.16). (Hint: first show inductively, using only the Peano axioms and basic set theory, that for every natural number $N \in \mathbf{N}$, there exists a unique pair A_N, B_N of subsets of \mathbf{N} which obeys the following properties: (a) $A_N \cap B_N = \emptyset$, (b) $A_N \cup B_N = \mathbf{N}$, (c) $0 \in A_N$, (d) $N++ \in B_N$, (e) Whenever $n \in B_N$, we have $n++ \in B_N$. (f) Whenever $n \in A_N$ and $n \neq N$, we have $n++ \in A_N$. Once one obtains these sets, use A_N as a substitute for $\{n \in \mathbf{N} : n \leq N\}$ in the previous argument.)

Exercise 3.5.13. The purpose of this exercise is to show that there is essentially only one version of the natural number system in set theory (cf. the discussion in Remark 2.1.12). Suppose we have a set \mathbf{N}' of “alternative natural numbers”, an “alternative zero” $0'$, and an “alternative increment operation” which takes any alternative natural number $n' \in \mathbf{N}'$ and returns another alternative

natural number $n'++' \in \mathbf{N}'$, such that the Peano Axioms (Axioms 2.1-2.5) all hold with the natural numbers, zero, and increment replaced by their alternative counterparts. Show that there exists a bijection $f : \mathbf{N} \rightarrow \mathbf{N}'$ from the natural numbers to the alternative natural numbers such that $f(0) = 0'$, and such that for any $n \in \mathbf{N}$ and $n' \in \mathbf{N}'$, we have $f(n) = n'$ if and only if $f(n++) = n'++'$. (Hint: use Exercise 3.5.12.)

3.6 Cardinality of sets

In the previous chapter we defined the natural numbers axiomatically, assuming that they were equipped with a 0 and an increment operation, and assuming five axioms on these numbers. Philosophically, this is quite different from one of our main conceptualizations of natural numbers - that of *cardinality*, or measuring *how many* elements there are in a set. Indeed, the Peano axiom approach treats natural numbers more like *ordinals* than *cardinals*. (The cardinals are One, Two, Three, ..., and are used to count how many things there are in a set. The *ordinals* are First, Second, Third, ..., and are used to order a sequence of objects. There is a subtle difference between the two, especially when comparing infinite cardinals with infinite ordinals, but this is beyond the scope of this text). We paid a lot of attention to what number came *next* after a given number n - which is an operation which is quite natural for ordinals, but less so for cardinals - but did not address the issue of whether these numbers could be used to *count* sets. The purpose of this section is to correct this issue by noting that the natural numbers *can* be used to count the cardinality of sets, as long as the set is finite.

The first thing is to work out when two sets have the same size: it seems clear that the sets $\{1, 2, 3\}$ and $\{4, 5, 6\}$ have the same size, but that both have a different size from $\{8, 9\}$. One way to define this is to say that two sets have the same size if they have the same number of elements, but we have not yet defined what the “number of elements” in a set is. Besides, this runs into problems when a set is infinite.

The right way to define the concept of “two sets having the same size” is not immediately obvious, but can be worked out with some thought. One intuitive reason why the sets $\{1, 2, 3\}$ and $\{4, 5, 6\}$ have the same size is that one can match the elements of the first set with the elements in the second set in a one-to-one correspondence: $1 \leftrightarrow 4$, $2 \leftrightarrow 5$, $3 \leftrightarrow 6$. (Indeed, this is how we first learn to count a set: we correspond the set we are trying to count with another set, such as a set of fingers on your hand). We will use this intuitive understanding as our rigorous basis for “having the same size”.

Definition 3.6.1 (Equal cardinality). We say that two sets X and Y have *equal cardinality* iff there exists a bijection $f : X \rightarrow Y$ from X to Y .

Example 3.6.2. The sets $\{0, 1, 2\}$ and $\{3, 4, 5\}$ have equal cardinality, since we can find a bijection between the two sets. Note that we do not yet know whether $\{0, 1, 2\}$ and $\{3, 4\}$ have equal cardinality; we know that one of the functions f from $\{0, 1, 2\}$ to $\{3, 4\}$ is not a bijection, but we have not proven yet that there might still be some other bijection from one set to the other. (It turns out that they do not have equal cardinality, but we will prove this a little later). Note that this definition makes sense regardless of whether X is finite or infinite (in fact, we haven’t even defined what finite means yet).

Note that two sets having equal cardinality does not preclude one set containing the other. For instance, if X is the set of natural numbers and Y is the set of even natural numbers, then the map $f : X \rightarrow Y$ defined by $f(n) := 2n$ is a bijection from X to Y (why?), and so X and Y have equal cardinality, despite Y being a subset of X and seeming intuitively as if it should only have “half” of the elements of X .

The notion of having equal cardinality is an equivalence relation:

Proposition 3.6.3. *Let X, Y, Z be sets. Then X has equal cardinality with X . If X has equal cardinality with Y , then Y has*

equal cardinality with X . If X has equal cardinality with Y and Y has equal cardinality with Z , then X has equal cardinality with Z .

Proof. See Exercise 3.6.1. □

Let n be a natural number. Now we want to say when a set X has n elements. Certainly we want the set $\{i \in \mathbf{N} : 1 \leq i \leq n\} = \{1, 2, \dots, n\}$ to have n elements. (This is true even when $n = 0$; the set $\{i \in \mathbf{N} : 1 \leq i \leq 0\}$ is just the empty set). Using our notion of equal cardinality, we thus define:

Definition 3.6.4. Let n be a natural number. A set X is said to have *cardinality* n , iff it has equal cardinality with $\{i \in \mathbf{N} : 1 \leq i \leq n\}$. We also say that X has n elements iff it has cardinality n .

Remark 3.6.5. One can use the set $\{i \in \mathbf{N} : i < n\}$ instead of $\{i \in \mathbf{N} : 1 \leq i \leq n\}$, since these two sets clearly have equal cardinality (why? What is the bijection?).

Example 3.6.6. Let a, b, c, d be distinct objects. Then $\{a, b, c, d\}$ has the same cardinality as $\{i \in \mathbf{N} : i < 4\} = \{0, 1, 2, 3\}$ or $\{i \in \mathbf{N} : 1 \leq i \leq 4\} = \{1, 2, 3, 4\}$ and thus has cardinality 4. Similarly, the set $\{a\}$ has cardinality 1.

There might be one problem with this definition: a set might have two different cardinalities. But this is not possible:

Proposition 3.6.7 (Uniqueness of cardinality). *Let X be a set with some cardinality n . Then X cannot have any other cardinality, i.e., X cannot have cardinality m for any $m \neq n$.*

Before we prove this proposition, we need a lemma.

Lemma 3.6.8. *Suppose that $n \geq 1$, and X has cardinality n . Then X is non-empty, and if x is any element of X , then the set $X - \{x\}$ (i.e., X with the element x removed) has cardinality $n - 1$.*

Proof. If X is empty then it clearly cannot have the same cardinality as the non-empty set $\{i \in \mathbf{N} : 1 \leq i \leq n\}$, as there is no bijection from the empty set to a non-empty set (why?). Now let x be an element of X . Since X has the same cardinality as $\{i \in \mathbf{N} : 1 \leq i \leq N\}$, we thus have a bijection f from X to $\{i \in \mathbf{N} : 1 \leq i \leq n\}$. In particular, $f(x)$ is a natural number between 1 and n . Now define the function $g : X - \{x\}$ to $\{i \in \mathbf{N} : 1 \leq i \leq n-1\}$ by the following rule: for any $y \in X - \{x\}$, we define $g(y) := f(y)$ if $f(y) < f(x)$, and define $g(y) := f(y) - 1$ if $f(y) > f(x)$. (Note that $f(y)$ cannot equal $f(x)$ since $y \neq x$ and f is a bijection.) It is easy to check that this map is also a bijection (why?), and so $X - \{x\}$ has equal cardinality with $\{i \in \mathbf{N} : 1 \leq i \leq n-1\}$. In particular $X - \{x\}$ has cardinality $n-1$, as desired. \square

Now we prove the proposition.

Proof of Proposition 3.6.7. We induct on n . First suppose that $n = 0$. Then X must be empty, and so X cannot have any non-zero cardinality. Now suppose that the Proposition is already proven for some n ; we now prove it for $n++$. Let X have cardinality $n++$; and suppose that X also has some other cardinality $m \neq n++$. By Proposition 3.6.3, X is non-empty, and if x is any element of X , then $X - \{x\}$ has cardinality n and also has cardinality $m-1$, by Lemma 3.6.8. By induction hypothesis, this means that $n = m-1$, which implies that $m = n++$, contradiction. This closes the induction. \square

Thus, for instance, we now know, thanks to Propositions 3.6.3 and 3.6.7, that the sets $\{0, 1, 2\}$ and $\{3, 4\}$ do not have equal cardinality, since the first set has cardinality 3 and the second set has cardinality 2.

Definition 3.6.9 (Finite sets). A set is *finite* iff it has cardinality n for some natural number n ; otherwise, the set is called *infinite*. If X is a finite set, we use $\#(X)$ to denote the cardinality of X .

Example 3.6.10. The sets $\{0, 1, 2\}$ and $\{3, 4\}$ are finite, as is the empty set (0 is a natural number), and $\#\{0, 1, 2\} = 3$, $\#\{3, 4\} = 2$, and $\#\emptyset = 0$.

Now we give an example of an infinite set.

Theorem 3.6.11. *The set of natural numbers \mathbf{N} is infinite.*

Proof. Suppose for contradiction that the set of natural numbers \mathbf{N} was finite, so it had some cardinality $\#\mathbf{N} = n$. Then there is a bijection f from $\{i \in \mathbf{N} : 1 \leq i \leq n\}$ to \mathbf{N} . One can show that the sequence $f(1), f(2), \dots, f(n)$ is bounded, or more precisely that there exists a natural number M such that $f(i) \leq M$ for all $1 \leq i \leq n$ (Exercise 3.6.3). But then the natural number $M + 1$ is not equal to any of the $f(i)$, contradicting the hypothesis that f is a bijection. \square

Remark 3.6.12. One can also use similar arguments to show that any unbounded set is infinite; for instance the rationals \mathbf{Q} and the reals \mathbf{R} (which we will construct in later chapters) are infinite. However, it is possible for some sets to be “more” infinite than others; see Section 8.3.

Now we relate cardinality with the arithmetic of natural numbers.

Proposition 3.6.13 (Cardinal arithmetic).

- (a) *Let X be a finite set, and let x be an object which is not an element of X . Then $X \cup \{x\}$ is finite and $\#(X \cup \{x\}) = \#(X) + 1$.*
- (b) *Let X and Y be finite sets. Then $X \cup Y$ is finite and $\#(X \cup Y) \leq \#(X) + \#(Y)$. If in addition X and Y are disjoint (i.e., $X \cap Y = \emptyset$), then $\#(X \cup Y) = \#(X) + \#(Y)$.*
- (c) *Let X be a finite set, and let Y be a subset of X . Then Y is finite, and $\#(Y) \leq \#(X)$. If in addition $Y \neq X$ (i.e., Y is a proper subset of X), then we have $\#(Y) < \#(X)$.*

- (d) If X is a finite set, and $f : X \rightarrow Y$ is a function, then $f(X)$ is a finite set with $\#(f(X)) \leq \#(X)$. If in addition f is one-to-one, then $\#(f(X)) = \#(X)$.
- (e) Let X and Y be finite sets. Then Cartesian product $X \times Y$ is finite and $\#(X \times Y) = \#(X) \times \#(Y)$.
- (f) Let X and Y be finite sets. Then the set Y^X (defined in Axiom 3.10) is finite and $\#(Y^X) = \#(Y)^{\#(X)}$.

Proof. See Exercise 3.6.4. □

Remark 3.6.14. Proposition 3.6.13 suggests that there is another way to define the arithmetic operations of natural numbers; not defined recursively as in Definitions 2.2.1, 2.3.1, 2.3.11, but instead using the notions of union, Cartesian product, and power set. This is the basis of *cardinal arithmetic*, which is an alternative foundation to arithmetic than the Peano arithmetic we have developed here; we will not develop this arithmetic in this text, but we give some examples of how one would work with this arithmetic in Exercises 3.6.5, 3.6.6.

This concludes our discussion of finite sets. We shall discuss infinite sets in Chapter 8, once we have constructed a few more examples of infinite sets (such as the integers, rationals and reals).

Exercise 3.6.1. Prove Proposition 3.6.3.

Exercise 3.6.2. Show that a set X has cardinality 0 if and only if X is the empty set.

Exercise 3.6.3. Let n be a natural number, and let $f : \{i \in \mathbf{N} : 1 \leq i \leq n\} \rightarrow \mathbf{N}$ be a function. Show that there exists a natural number M such that $f(i) \leq M$ for all $1 \leq i \leq n$. (Hint: induct on n . You may also want to peek at Lemma 5.1.14.) This exercise shows that finite subsets of the natural numbers are bounded.

Exercise 3.6.4. Prove Proposition 3.6.13.

Exercise 3.6.5. Let A and B be sets. Show that $A \times B$ and $B \times A$ have equal cardinality by constructing an explicit bijection between the two sets. Then use Proposition 3.6.13 to conclude an

alternate proof of Lemma 2.3.2. (We remark that many other results in that section could also be proven by similar methods.)

Exercise 3.6.6. Let A, B, C be sets. Show that the sets $(A^B)^C$ and $A^{B \times C}$ have equal cardinality by constructing an explicit bijection between the two sets. If B and C are disjoint, show that $A^B \times A^C$ and $A^{B \cup C}$ also have equal cardinality. Then use Proposition 3.6.13 to conclude that for any natural numbers a, b, c , that $(a^b)^c = a^{bc}$ and $a^b \times a^c = a^{b+c}$.

Exercise 3.6.7. Let A and B be sets. Let us say that A has *lesser or equal* cardinality to B if there exists an injection $f : A \rightarrow B$ from A to B . Show that if A and B are finite sets, then A has lesser or equal cardinality to B if and only if $\#(A) \leq \#(B)$.

Exercise 3.6.8. Let A and B be sets such that there exists an injection $f : A \rightarrow B$ from A to B (i.e., A has lesser or equal cardinality to B). Show that there then exists a surjection $g : B \rightarrow A$ from B to A . (The converse to this statement requires the axiom of choice; see Exercise 8.4.3.)

Exercise 3.6.9. Let A and B be finite sets. Show that $A \cup B$ and $A \cap B$ are also finite sets, and that $\#(A) + \#(B) = \#(A \cup B) + \#(A \cap B)$.

Chapter 4

Integers and rationals

4.1 The integers

In Chapter 2 we built up most of the basic properties of the natural number system, but are reaching the limits of what one can do with just addition and multiplication. We would now like to introduce a new operation, that of subtraction, but to do that properly we will have to pass from the natural number system to a larger number system, that of the *integers*.

Informally, the integers are what you can get by subtracting two natural numbers; for instance, $3 - 5$ should be an integer, as should $6 - 2$. This is not a complete definition of the integers, because (a) it doesn't say when two differences are equal (for instance we should know why $3 - 5$ is equal to $2 - 4$, but is not equal to $1 - 6$), and (b) it doesn't say how to do arithmetic on these differences (how does one add $3 - 5$ to $6 - 2$?). Furthermore, (c) this definition is circular because it requires a notion of subtraction, which we can only adequately define once the integers are constructed. Fortunately, because of our prior experience with integers we know what the answers to these questions should be. To answer (a), we know from our advanced knowledge in algebra that $a - b = c - d$ happens exactly when $a + d = c + b$, so we can characterize equality of differences using only the concept of addition. Similarly, to answer (b) we know from algebra that $(a - b) + (c - d) = (a + c) - (b + d)$ and that $(a - b)(c - d) = (ac + bd) - (ad + bc)$. So we will take advan-

tage of our foreknowledge by building all this into the *definition* of the integers, as we shall do shortly.

We still have to resolve (c). To get around this problem we will use the following work-around: we will temporarily write integers not as a difference $a - b$, but instead use a new notation $a \text{---} b$ to define integers, where the --- is a meaningless place-holder (similar to the comma in the Cartesian co-ordinate notation (x, y) for points in the plane). Later when we define subtraction we will see that $a \text{---} b$ is in fact equal to $a - b$, and so we can discard the notation --- ; it is only needed right now to avoid circularity. (These devices are similar to the scaffolding used to construct a building; they are temporarily essential to make sure the building is built correctly, but once the building is completed they are thrown away and never used again). This may seem unnecessarily complicated in order to define something that we already are very familiar with, but we will use this device again to construct the rationals, and knowing these kinds of constructions will be very helpful in later chapters.

Definition 4.1.1 (Integers). An *integer* is an expression¹ of the form $a \text{---} b$, where a and b are natural numbers. Two integers are considered to be equal, $a \text{---} b = c \text{---} d$, if and only if $a + d = c + b$. We let \mathbf{Z} denote the set of all integers.

Thus for instance $3 \text{---} 5$ is an integer, and is equal to $2 \text{---} 4$, because $3 + 4 = 2 + 5$. On the other hand, $3 \text{---} 5$ is not equal to $2 \text{---} 3$ because $3 + 3 \neq 2 + 5$. (This notation is strange looking, and has a few deficiencies; for instance, 3 is not yet an integer, because it is not of the form $a \text{---} b$! We will rectify these problems later.)

¹In the language of set theory, what we are doing here is starting with the space $\mathbf{N} \times \mathbf{N}$ of ordered pairs (a, b) of natural numbers. Then we place an equivalence relation \sim on these pairs by declaring $(a, b) \sim (c, d)$ iff $a + d = c + b$. The set-theoretic interpretation of the symbol $a \text{---} b$ is that it is the space of all pairs equivalent to (a, b) : $a \text{---} b := \{(c, d) \in \mathbf{N} \times \mathbf{N} : (a, b) \sim (c, d)\}$. However, this interpretation plays no role on how we manipulate the integers and we will not refer to it again. A similar set-theoretic interpretation can be given to the construction of the rational numbers later in this chapter, or the real numbers in the next chapter.

We have to check that this is a legitimate notion of equality. We need to verify the reflexivity, symmetry, transitivity, and substitution axioms (see Section 12.7). We leave reflexivity and symmetry to Exercise 4.1.1 and instead verify the transitivity axiom. Suppose we know that $a \text{---} b = c \text{---} d$ and $c \text{---} d = e \text{---} f$. Then we have $a + d = c + b$ and $c + f = d + e$. Adding the two equations together we obtain $a + d + c + f = c + b + d + e$. By Proposition 2.2.6 we can cancel the c and d , obtaining $a + f = b + e$, i.e., $a \text{---} b = e \text{---} f$. Thus the cancellation law was needed to make sure that our notion of equality is sound. As for the substitution axiom, we cannot verify it at this stage because we have not yet defined any operations on the integers. However, when we do define our basic operations on the integers, such as addition, multiplication, and order, we will have to verify the substitution axiom at that time in order to ensure that the definition is valid. (We will only need to do this for the basic operations; more advanced operations on the integers, such as exponentiation, will be defined in terms of the basic ones, and so we do not need to re-verify the substitution axiom for the advanced operations.)

Now we define two basic arithmetic operations on integers: addition and multiplication.

Definition 4.1.2. The sum of two integers, $(a \text{---} b) + (c \text{---} d)$, is defined by the formula

$$(a \text{---} b) + (c \text{---} d) := (a + c) \text{---} (b + d).$$

The product of two integers, $(a \text{---} b) \times (c \text{---} d)$, is defined by

$$(a \text{---} b) \times (c \text{---} d) := (ac + bd) \text{---} (ad + bc).$$

Thus for instance, $(3 \text{---} 5) + (1 \text{---} 4)$ is equal to $(4 \text{---} 9)$. There is however one thing we have to check before we can accept these definitions - we have to check that if we replace one of the integers by an equal integer, that the sum or product does not change. For instance, $(3 \text{---} 5)$ is equal to $(2 \text{---} 4)$, so $(3 \text{---} 5) + (1 \text{---} 4)$ ought to have the same value as $(2 \text{---} 4) + (1 \text{---} 4)$, otherwise this would not give a consistent definition of addition. Fortunately, this is the case:

Lemma 4.1.3 (Addition and multiplication are well-defined). *Let a, b, a', b', c, d be natural numbers. If $(a—b) = (a'—b')$, then $(a—b) + (c—d) = (a'—b') + (c—d)$ and $(a—b) \times (c—d) = (a'—b') \times (c—d)$, and also $(c—d) + (a—b) = (c—d) + (a'—b')$ and $(c—d) \times (a—b) = (c—d) \times (a'—b')$. Thus addition and multiplication are well-defined operations (equal inputs give equal outputs).*

Proof. To prove that $(a—b) + (c—d) = (a'—b') + (c—d)$, we evaluate both sides as $(a + c)—(b + d)$ and $(a' + c)—(b' + d)$. Thus we need to show that $a + c + b' + d = a' + c + b + d$. But since $(a—b) = (a'—b')$, we have $a + b' = a' + b$, and so by adding $c + d$ to both sides we obtain the claim. Now we show that $(a—b) \times (c—d) = (a'—b') \times (c—d)$. Both sides evaluate to $(ac + bd)—(ad + bc)$ and $(a'c + b'd)—(a'd + b'c)$, so we have to show that $ac + bd + a'd + b'c = a'c + b'd + ad + bc$. But the left-hand side factors as $c(a + b') + d(a' + b)$, while the right factors as $c(a' + b) + d(a + b')$. Since $a + b' = a' + b$, the two sides are equal. The other two identities are proven similarly. \square

The integers $n—0$ behave in the same way as the natural numbers n ; indeed one can check that $(n—0) + (m—0) = (n + m)—0$ and $(n—0) \times (m—0) = nm—0$. Furthermore, $(n—0)$ is equal to $(m—0)$ if and only if $n = m$. (The mathematical term for this is that there is an *isomorphism* between the natural numbers n and those integers of the form $n—0$). Thus we may *identify* the natural numbers with integers by setting $n \equiv n—0$; this does not affect our definitions of addition or multiplication or equality since they are consistent with each other. Thus for instance the natural number 3 is now considered to be the same as the integer $3—0$: $3 = 3—0$. In particular 0 is equal to $0—0$ and 1 is equal to $1—0$. Of course, if we set n equal to $n—0$, then it will also be equal to any other integer which is equal to $n—0$, for instance 3 is equal not only to $3—0$, but also to $4—1$, $5—2$, etc.

We can now define incrementation on the integers by defining $x++ := x + 1$ for any integer x ; this is of course consistent with

our definition of the increment operation for natural numbers. However, this is no longer an important operation for us, as it has been now superseded by the more general notion of addition.

Now we consider some other basic operations on the integers.

Definition 4.1.4 (Negation of integers). If $(a—b)$ is an integer, we define the negation $-(a—b)$ to be the integer $(b—a)$. In particular if $n = n—0$ is a positive natural number, we can define its negation $-n = 0—n$.

For instance $-(3—5) = (5—3)$. One can check this definition is well-defined (Exercise 4.1.2).

We can now show that the integers correspond exactly to what we expect.

Lemma 4.1.5 (Trichotomy of integers). *Let x be an integer. Then exactly one of the following three statements is true: (a) x is zero; (b) x is equal to a positive natural number n ; or (c) x is the negation $-n$ of a positive natural number n .*

Proof. We first show that at least one of (a), (b), (c) is true. By definition, $x = a—b$ for some natural numbers a, b . We have three cases: $a > b$, $a = b$, or $a < b$. If $a > b$ then $a = b + c$ for some positive natural number c , which means that $a—b = c—0 = c$, which is (b). If $a = b$, then $a—b = a—a = 0—0 = 0$, which is (a). If $a < b$, then $b > a$, so that $b—a = n$ for some natural number n by the previous reasoning, and thus $a—b = -n$, which is (c).

Now we show that no more than one of (a), (b), (c) can hold at a time. By definition, a positive natural number is non-zero, so (a) and (b) cannot simultaneously be true. If (a) and (c) were simultaneously true, then $0 = -n$ for some positive natural n ; thus $(0—0) = (0—n)$, so that $0 + n = 0 + 0$, so that $n = 0$, a contradiction. If (b) and (c) were simultaneously true, then $n = -m$ for some positive n, m , so that $(n—0) = (0—m)$, so that $n + m = 0 + 0$, which contradicts Proposition 2.2.8. Thus exactly one of (a), (b), (c) is true for any integer x . \square

If n is a positive natural number, we call $-n$ a *negative integer*. Thus every integer is positive, zero, or negative, but not more than one of these at a time.

One could well ask why we don't use Lemma 4.1.5 to *define* the integers; i.e., why didn't we just say an integer is anything which is either a positive natural number, zero, or the negative of a natural number. The reason is that if we did so, the rules for adding and multiplying integers would split into many different cases (e.g., negative times positive equals positive; negative plus positive is either negative, positive, or zero, depending on which term is larger, etc.) and to verify all the properties ends up being much messier than doing it this way.

We now summarize the algebraic properties of the integers.

Proposition 4.1.6 (Laws of algebra for integers). *Let x, y, z be integers. Then we have*

$$\begin{aligned}x + y &= y + x \\(x + y) + z &= x + (y + z) \\x + 0 &= 0 + x = x \\x + (-x) &= (-x) + x = 0 \\xy &= yx \\(xy)z &= x(yz) \\x1 &= 1x = x \\x(y + z) &= xy + xz \\(y + z)x &= yx + zx.\end{aligned}$$

Remark 4.1.7. The above set of nine identities have a name; they are asserting that the integers form a *commutative ring*. (If one deleted the identity $xy = yx$, then they would only assert that the integers form a *ring*). Note that some of these identities were already proven for the natural numbers, but this does not automatically mean that they also hold for the integers because the integers are a larger set than the natural numbers. On the other hand, this Proposition supercedes many of the propositions derived earlier for natural numbers.

Proof. There are two ways to prove these identities. One is to use Lemma 4.1.5 and split into a lot of cases depending on whether x, y, z are zero, positive, or negative. This becomes very messy. A shorter way is to write $x = (a \text{---} b)$, $y = (c \text{---} d)$, and $z = (e \text{---} f)$ for some natural numbers a, b, c, d, e, f , and expand these identities in terms of a, b, c, d, e, f and use the algebra of the natural numbers. This allows each identity to be proven in a few lines. We shall just prove the longest one, namely $(xy)z = x(yz)$:

$$\begin{aligned} (xy)z &= ((a \text{---} b)(c \text{---} d))(e \text{---} f) \\ &= ((ac + bd) \text{---} (ad + bc))(e \text{---} f) \\ &= ((ace + bde + adf + bcf) \text{---} (acf + bdf + ade + bce)); \\ x(yz) &= (a \text{---} b)((c \text{---} d)(e \text{---} f)) \\ &= (a \text{---} b)((ce + df) \text{---} (cf + de)) \\ &= ((ace + adf + bcf + bde) \text{---} (acf + ade + bcd + bdf)) \end{aligned}$$

and so one can see that $(xy)z$ and $x(yz)$ are equal. The other identities are proven in a similar fashion; see Exercise 4.1.4. \square

We now define the operation of *subtraction* $x - y$ of two integers by the formula

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

We do not need to verify the substitution axiom for this operation, since we have defined subtraction in terms of two other operations on integers, namely addition and negation, and we have already verified that those operations are well-defined.

One can easily check now that if a and b are natural numbers, then

$$a - b = a + -b = (a \text{---} 0) + (0 \text{---} b) = a \text{---} b,$$

and so $a \text{---} b$ is just the same thing as $a - b$. Because of this we can now discard the --- notation, and use the familiar operation of subtraction instead. (As remarked before, we could not use subtraction immediately because it would be circular.)

We can now generalize Lemma 2.3.3 and Corollary 2.3.7 from the natural numbers to the integers:

Proposition 4.1.8 (Integers have no zero divisors). *Let a and b be integers such that $ab = 0$. Then either $a = 0$ or $b = 0$ (or both).*

Proof. See Exercise 4.1.5. □

Corollary 4.1.9 (Cancellation law for integers). *If a, b, c are integers such that $ac = bc$ and c is non-zero, then $a = b$.*

Proof. See Exercise 4.1.6. □

We now extend the notion of order, which was defined on the natural numbers, to the integers by repeating the definition verbatim:

Definition 4.1.10 (Ordering of the integers). Let n and m be integers. We say that n is *greater than or equal to* m , and write $n \geq m$ or $m \leq n$, iff we have $n = m + a$ for some natural number a . We say that n is *strictly greater than* m , and write $n > m$ or $m < n$, iff $n \geq m$ and $n \neq m$.

Thus for instance $5 > -3$, because $5 = -3 + 8$ and $5 \neq -3$. Clearly this definition is consistent with the notion of order on the natural numbers, since we are using the same definition.

Using the laws of algebra in Proposition 4.1.6 it is not hard to show the following properties of order:

Lemma 4.1.11 (Properties of order). *Let a, b, c be integers.*

- $a > b$ if and only if $a - b$ is a positive natural number.
- (Addition preserves order) If $a > b$, then $a + c > b + c$.
- (Positive multiplication preserves order) If $a > b$ and c is positive, then $ac > bc$.
- (Negation reverses order) If $a > b$, then $-a < -b$.
- (Order is transitive) If $a > b$ and $b > c$, then $a > c$.
- (Order trichotomy) Exactly one of the statements $a > b$, $a < b$, or $a = b$ is true.

Proof. See Exercise 4.1.7. □

Exercise 4.1.1. Verify that the definition of equality on the integers is both reflexive and symmetric.

Exercise 4.1.2. Show that the definition of negation on the integers is well-defined in the sense that if $(a - b) = (a' - b')$, then $-(a - b) = -(a' - b')$ (so equal integers have equal negations).

Exercise 4.1.3. Show that $(-1) \times a = -a$ for every integer a .

Exercise 4.1.4. Prove the remaining identities in Proposition 4.1.6. (Hint: one can save some work by using some identities to prove others. For instance, once you know that $xy = yx$, you get for free that $x1 = 1x$, and once you also prove $x(y + z) = xy + xz$, you automatically get $(y + z)x = yx + zx$ for free.)

Exercise 4.1.5. Prove Proposition 4.1.8. (Hint: while this Proposition is not quite the same as Lemma 2.3.3, it is certainly legitimate to use Lemma 2.3.3 in the course of proving Proposition 4.1.8.)

Exercise 4.1.6. Prove Corollary 4.1.9. (Hint: There are two ways to do this. One is to use Proposition 4.1.8 to conclude that $a - b$ must be zero. Another way is to combine Corollary 2.3.7 with Lemma 4.1.5.)

Exercise 4.1.7. Prove Lemma 4.1.11. (Hint: use the first part of this Lemma to prove all the others.)

Exercise 4.1.8. Show that the principle of induction (Axiom V) does not apply directly to the integers. More precisely, give an example of a property $P(n)$ pertaining to an integer n such that $P(0)$ is true, and that $P(n)$ implies $P(n++)$ for all integers n , but that $P(n)$ is not true for all integers n . Thus induction is not as useful a tool for dealing with the integers as it is with the natural numbers. (The situation becomes even worse with the rational and real numbers, which we shall define shortly.)

4.2 The rationals

We have now constructed the integers, with the operations of addition, subtraction, multiplication, and order and verified all the

expected algebraic and order-theoretic properties. Now we will make a similar construction to build the rationals, adding division to our mix of operations.

Just like the integers were constructed by subtracting two natural numbers, the rationals can be constructed by dividing two integers, though of course we have to make the usual caveat that the denominator should² be non-zero. Of course, just as two differences $a - b$ and $c - d$ can be equal if $a + d = c + b$, we know (from more advanced knowledge) that two quotients a/b and c/d can be equal if $ad = bc$. Thus, in analogy with the integers, we create a new meaningless symbol $//$ (which will eventually be superceded by division), and define

Definition 4.2.1. A *rational number* is an expression of the form $a//b$, where a and b are integers and b is non-zero; $a//0$ is not considered to be a rational number. Two rational numbers are considered to be equal, $a//b = c//d$, if and only if $ad = cb$. The set of all rational numbers is denoted \mathbf{Q} .

Thus for instance $3//4 = 6//8 = -3// -4$, but $3//4 \neq 4//3$. This is a valid definition of equality (Exercise 4.2.1). Now we need a notion of addition, multiplication, and negation. Again, we will take advantage of our pre-existing knowledge, which tells us that $a/b + c/d$ should equal $(ad + bc)/(bd)$ and that $a/b * c/d$ should equal ac/bd , while $-(a/b)$ equals $(-a)/b$. Motivated by this foreknowledge, we define

Definition 4.2.2. If $a//b$ and $c//d$ are rational numbers, we define their sum

$$(a//b) + (c//d) := (ad + bc)//(bd)$$

their product

$$(a//b) * (c//d) := (ac)//(bd)$$

²There is no reasonable way we can divide by zero, since one cannot have both the identities $(a/b)*b = a$ and $c*0 = 0$ hold simultaneously if b is allowed to be zero. However, we can eventually get a reasonable notion of dividing by a quantity which *approaches* zero - think of L'Hôpital's rule (see Section 10.4), which suffices for doing things like defining differentiation.

and the negation

$$-(a//b) := (-a)//b.$$

Note that if b and d are non-zero, then bd is also non-zero, by Proposition 4.1.8, so the sum or product of a rational number remains a rational number.

Lemma 4.2.3. *The sum, product, and negation operations on rational numbers are well-defined, in the sense that if one replaces $a//b$ with another rational number $a'//b'$ which is equal to $a//b$, then the output of the above operations remains unchanged, and similarly for $c//d$.*

Proof. We just verify this for addition; we leave the remaining claims to Exercise 4.2.2. Suppose $a//b = a'//b'$, so that b and b' are non-zero and $ab' = a'b$. We now show that $a//b + c//d = a'//b' + c//d$. By definition, the left-hand side is $(ad+bc)//bd$ and the right-hand side is $(a'd + b'c)//b'd$, so we have to show that

$$(ad + bc)b'd = (a'd + b'c)bd,$$

which expands to

$$ab'd^2 + bb'cd = a'bd^2 + bb'cd.$$

But since $ab' = a'b$, the claim follows. Similarly if one replaces $c//d$ by $c'//d'$. \square

We note that the rational numbers $a//1$ behave in a manner identical to the integers a :

$$\begin{aligned}(a//1) + (b//1) &= (a + b)//1; \\ (a//1) \times (b//1) &= (ab)//1; \\ -(a//1) &= (-a)//1.\end{aligned}$$

Also, $a//1$ and $b//1$ are only equal when a and b are equal. Because of this, we will identify a with $a//1$ for each integer a : $a \equiv a//1$; the above identities then guarantee that the arithmetic of the integers is consistent with the arithmetic of the rationals.

Thus just as we embedded the natural numbers inside the integers, we embed the integers inside the rational numbers. In particular, all natural numbers are rational numbers, for instance 0 is equal to $0//1$ and 1 is equal to $1//1$.

Observe that a rational number $a//b$ is equal to $0 = 0//1$ if and only if $a \times 1 = b \times 0$, i.e., if the numerator a is equal to 0. Thus if a and b are non-zero then so is $a//b$.

We now define a new operation on the rationals: reciprocal. If $x = a//b$ is a non-zero rational (so that $a, b \neq 0$) then we define the *reciprocal* x^{-1} of x to be the rational number $x^{-1} := b//a$. It is easy to check that this operation is consistent with our notion of equality: if two rational numbers $a//b, a'//b'$ are equal, then their reciprocals are also equal. (In contrast, an operation such as “numerator” is not well-defined: the rationals $3//4$ and $6//8$ are equal, but have unequal numerators, so we have to be careful when referring to such terms as “the numerator of x ”.) We however leave the reciprocal of 0 undefined.

We now summarize the algebraic properties of the rationals.

Proposition 4.2.4 (Laws of algebra for rationals). *Let x, y, z be rationals. Then the following laws of algebra hold:*

$$\begin{aligned} x + y &= y + x \\ (x + y) + z &= x + (y + z) \\ x + 0 &= 0 + x = x \\ x + (-x) &= (-x) + x = 0 \\ xy &= yx \\ (xy)z &= x(yz) \\ x1 &= 1x = x \\ x(y + z) &= xy + xz \\ (y + z)x &= yx + zx. \end{aligned}$$

If x is non-zero, we also have

$$xx^{-1} = x^{-1}x = 1.$$

Remark 4.2.5. The above set of ten identities have a name; they are asserting that the rationals \mathbf{Q} form a *field*. This is better than being a commutative ring because of the tenth identity $xx^{-1} = x^{-1}x = 1$. Note that this Proposition supercedes Proposition 4.1.6.

Proof. To prove this identity, one writes $x = a//b$, $y = c//d$, $z = e//f$ for some integers a, c, e and non-zero integers b, d, f , and verifies each identity in turn using the algebra of the integers. We shall just prove the longest one, namely $(x+y)+z = x+(y+z)$:

$$\begin{aligned} (x+y)+z &= ((a//b) + (c//d)) + (e//f) \\ &= ((ad+bc)//bd) + (e//f) \\ &= (adf+bcf+bde)//bdf; \\ x+(y+z) &= (a//b) + ((c//d) + (e//f)) \\ &= (a//b) + ((cf+de)//df) = (adf+bcf+bde)//bdf \end{aligned}$$

and so one can see that $(x+y)+z$ and $x+(y+z)$ are equal. The other identities are proven in a similar fashion and are left to Exercise 4.2.4. \square

We can now define the *quotient* x/y of two rational numbers x and y , provided that y is non-zero, by the formula

$$x/y := x \times y^{-1}.$$

Thus, for instance

$$(3//4)/(5//6) = (3//4) \times (6//5) = (18//20) = (9//10).$$

Using this formula, it is easy to see that $a/b = a//b$ for every integer a and every non-zero integer b . Thus we can now discard the $//$ notation, and use the more customary a/b instead of $a//b$.

The above field axioms allow us to use all the normal rules of algebra; we will now proceed to do so without further comment.

In the previous section we organized the integers into positive, zero, and negative numbers. We now do the same for the rationals.

Definition 4.2.6. A rational number x is said to be *positive* iff we have $x = a/b$ for some positive integers a and b . It is said to be *negative* iff we have $x = -y$ for some positive rational y (i.e., $x = (-a)/b$ for some positive integers a and b).

Thus for instance, every positive integer is a positive rational number, and every negative integer is a negative rational number, so our new definition is consistent with our old one.

Lemma 4.2.7 (Trichotomy of rationals). *Let x be a rational number. Then exactly one of the following three statements is true: (a) x is equal to 0. (b) x is a positive rational number. (c) x is a negative rational number.*

Proof. See Exercise 4.2.4. □

Definition 4.2.8 (Ordering of the rationals). Let x and y be rational numbers. We say that $x > y$ iff $x - y$ is a positive rational number, and $x < y$ iff $x - y$ is a negative rational number. We write $x \geq y$ iff either $x > y$ or $x = y$, and similarly define $x \leq y$.

Proposition 4.2.9 (Basic properties of order on the rationals). *Let x, y, z be rational numbers. Then the following properties hold.*

- (a) (Order trichotomy) *Exactly one of the three statements $x = y$, $x < y$, or $x > y$ is true.*
- (b) (Order is anti-symmetric) *One has $x < y$ if and only if $y > x$.*
- (c) (Order is transitive) *If $x < y$ and $y < z$, then $x < z$.*
- (d) (Addition preserves order) *If $x < y$, then $x + z < y + z$.*
- (e) (Positive multiplication preserves order) *If $x < y$ and z is positive, then $xz < yz$.*

Proof. See Exercise 4.2.5. □

Remark 4.2.10. The above five properties in Proposition 4.2.9, combined with the field axioms in Proposition 4.2.4, have a name: they assert that the rationals \mathbf{Q} form an *ordered field*. It is important to keep in mind that Proposition 4.2.9(e) only works when z is positive, see Exercise 4.2.6.

Exercise 4.2.1. Show that the definition of equality for the rational numbers is reflexive, symmetric, and transitive. (Hint: for transitivity, use Corollary 2.3.7.)

Exercise 4.2.2. Prove the remaining components of Lemma 4.2.3.

Exercise 4.2.3. Prove the remaining components of Proposition 4.2.4. (Hint: as with Proposition 4.1.6, you can save some work by using some identities to prove others.)

Exercise 4.2.4. Prove Lemma 4.2.7. (Note that, as in Proposition 2.2.13, you have to prove two different things: firstly, that *at least* one of (a), (b), (c) is true; and secondly, that *at most* one of (a), (b), (c) is true.)

Exercise 4.2.5. Prove Proposition 4.2.9.

Exercise 4.2.6. Show that if x, y, z are real numbers such that $x < y$ and z is *negative*, then $xz > yz$.

4.3 Absolute value and exponentiation

We have already introduced the four basic arithmetic operations of addition, subtraction, multiplication, and division on the rationals. (Recall that subtraction and division came from the more primitive notions of negation and reciprocal by the formulae $x - y := x + (-y)$ and $x/y := x \times y^{-1}$.) We also have a notion of order $<$, and have organized the rationals into the positive rationals, the negative rationals, and zero. In short, we have shown that the rationals \mathbf{Q} form an *ordered field*.

One can now use these basic operations to construct more operations. There are many such operations we can construct, but we shall just introduce two particularly useful ones: absolute value and exponentiation.

Definition 4.3.1 (Absolute value). If x is a rational number, the *absolute value* $|x|$ of x is defined as follows. If x is positive, then $|x| := x$. If x is negative, then $|x| := -x$. If x is zero, then $|x| := 0$.

Definition 4.3.2 (Distance). Let x and y be real numbers. The quantity $|x - y|$ is called the *distance between x and y* and is sometimes denoted $d(x, y)$, thus $d(x, y) := |x - y|$. For instance, $d(3, 5) = 2$.

Proposition 4.3.3 (Basic properties of absolute value and distance). *Let x, y, z be rational numbers.*

- (a) (Non-degeneracy of absolute value) *We have $|x| \geq 0$. Also, $|x| = 0$ if and only if x is 0.*
- (b) (Triangle inequality for absolute value) *We have $|x + y| \leq |x| + |y|$.*
- (c) *We have the inequalities $-y \leq x \leq y$ if and only if $y \geq |x|$. In particular, we have $-|x| \leq x \leq |x|$.*
- (d) (Multiplicativity of absolute value) *We have $|xy| = |x| |y|$. In particular, $|-x| = |x|$.*
- (e) (Non-degeneracy of distance) *We have $d(x, y) \geq 0$. Also, $d(x, y) = 0$ if and only if $x = y$.*
- (f) (Symmetry of distance) $d(x, y) = d(y, x)$
- (g) (Triangle inequality for distance) $d(x, z) \leq d(x, y) + d(y, z)$.

Proof. See Exercise 4.3.1. □

Absolute value is useful for measuring how “close” two numbers are. Let us make a somewhat artificial definition:

Definition 4.3.4 (ε -closeness). Let $\varepsilon > 0$, and x, y be rational numbers. We say that y is ε -close to x iff we have $d(y, x) \leq \varepsilon$.

Remark 4.3.5. This definition is not standard in mathematics textbooks; we will use it as “scaffolding” to construct the more important notions of limits (and of Cauchy sequences) later on, and once we have those more advanced notions we will discard the notion of ε -close.

Examples 4.3.6. The numbers 0.99 and 1.01 are 0.1-close, but they are not 0.01 close, because $d(0.99, 1.01) = |0.99 - 1.01| = 0.02$ is larger than 0.01. The numbers 2 and 2 are ε -close for every positive ε .

We do not bother defining a notion of ε -close when ε is zero or negative, because if ε is zero then x and y are only ε -close when they are equal, and when ε is negative then x and y are never ε -close. (In any event it is a long-standing tradition in analysis that the Greek letters ε , δ should only denote positive (and probably small) numbers).

Some basic properties of ε -closeness are the following.

Proposition 4.3.7. *Let x, y, z, w be rational numbers.*

- (a) *If $x = y$, then x is ε -close to y for every $\varepsilon > 0$. Conversely, if x is ε -close to y for every $\varepsilon > 0$, then we have $x = y$.*
- (b) *Let $\varepsilon > 0$. If x is ε -close to y , then y is ε -close to x .*
- (c) *Let $\varepsilon, \delta > 0$. If x is ε -close to y , and y is δ -close to z , then x and z are $(\varepsilon + \delta)$ -close.*
- (d) *Let $\varepsilon, \delta > 0$. If x and y are ε -close, and z and w are δ -close, then $x + z$ and $y + w$ are $(\varepsilon + \delta)$ -close, and $x - z$ and $y - w$ are also $(\varepsilon + \delta)$ -close.*
- (e) *Let $\varepsilon > 0$. If x and y are ε -close, they are also ε' -close for every $\varepsilon' > \varepsilon$.*
- (f) *Let $\varepsilon > 0$. If y and z are both ε -close to x , and w is between y and z (i.e., $y \leq w \leq z$ or $z \leq w \leq y$), then w is also ε -close to x .*

(g) Let $\varepsilon > 0$. If x and y are ε -close, and z is non-zero, then xz and yz are $\varepsilon|z|$ -close.

(h) Let $\varepsilon, \delta > 0$. If x and y are ε -close, and z and w are δ -close, then xz and yw are $(\varepsilon|z| + \delta|x| + \varepsilon\delta)$ -close.

Proof. We only prove the most difficult one, (h); we leave (a)-(g) to Exercise 4.3.2. Let $\varepsilon, \delta > 0$, and suppose that x and y are ε -close. If we write $a := y - x$, then we have $y = x + a$ and that $|a| \leq \varepsilon$. Similarly, if z and w are δ -close, and we define $b := w - z$, then $w = z + b$ and $|b| \leq \delta$.

Since $y = x + a$ and $w = z + b$, we have

$$yw = (x + a)(z + b) = xz + az + xb + ab.$$

Thus

$$|yw - xz| = |az + bx + ab| \leq |az| + |bx| + |ab| = |a||z| + |b||x| + |a||b|.$$

Since $|a| \leq \varepsilon$ and $|b| \leq \delta$, we thus have

$$|yw - xz| \leq \varepsilon|z| + \delta|x| + \varepsilon\delta$$

and thus that yw and xz are $(\varepsilon|z| + \delta|x| + \varepsilon\delta)$ -close. \square

Remark 4.3.8. One should compare statements (a)-(c) of this Proposition with the reflexive, symmetric, and transitive axioms of equality. It is often useful to think of the notion of “ ε -close” as an approximate substitute for that of equality in analysis.

Now we recursively define exponentiation for natural number exponents, extending the previous definition in Definition 2.3.11.

Definition 4.3.9 (Exponentiation to a natural number). Let x be a rational number. To raise x to the power 0, we define $x^0 := 1$. Now suppose inductively that x^n has been defined for some natural number n , then we define $x^{n+1} := x^n \times x$.

Proposition 4.3.10 (Properties of exponentiation, I). Let x, y be rational numbers, and let n, m be natural numbers.

- (a) We have $x^n x^m = x^{n+m}$, $(x^n)^m = x^{nm}$, and $(xy)^n = x^n y^n$.
- (b) We have $x^n = 0$ if and only if $x = 0$.
- (c) If $x \geq y \geq 0$, then $x^n \geq y^n \geq 0$. If $x > y \geq 0$, then $x^n > y^n \geq 0$.
- (d) We have $|x^n| = |x|^n$.

Proof. See Exercise 4.3.3. □

Now we define exponentiation for negative integer exponents.

Definition 4.3.11 (Exponentiation to a negative number). Let x be a non-zero rational number. Then for any negative integer $-n$, we define $x^{-n} := 1/x^n$.

Thus for instance $x^{-3} = 1/x^3 = 1/(x \times x \times x)$. We now have x^n defined for any integer n , whether n is positive, negative, or zero. Exponentiation with integer exponents has the following properties (which supercede Proposition 4.3.10):

Proposition 4.3.12 (Properties of exponentiation, II). *Let x, y be non-zero rational numbers, and let n, m be integers.*

- (a) We have $x^n x^m = x^{n+m}$, $(x^n)^m = x^{nm}$, and $(xy)^n = x^n y^n$.
- (b) If $x \geq y > 0$, then $x^n \geq y^n > 0$ if n is positive, and $0 < x^n \leq y^n$ if n is negative.
- (c) If $x, y \geq 0$, $n \neq 0$, and $x^n = y^n$, then $x = y$.
- (d) We have $|x^n| = |x|^n$.

Proof. See Exercise 4.3.4. □

Exercise 4.3.1. Prove Proposition 4.3.3. (Hint: While all of these claims can be proven by dividing into cases, such as when x is positive, negative, or zero, several parts of the proposition can be proven without such tedious dividing into cases, for instance by using earlier parts of the proposition to prove later ones.)

Exercise 4.3.2. Prove the remaining claims in Proposition 4.3.7.

Exercise 4.3.3. Prove Proposition 4.3.10. (Hint: use induction.)

Exercise 4.3.4. Prove Proposition 4.3.12. (Hint: induction is not suitable here. Instead, use Proposition 4.3.10).

Exercise 4.3.5. Prove that $2^N \geq N$ for all positive integers N . (Hint: use induction).

4.4 Gaps in the rational numbers

Imagine that we arrange the rationals on a line, arranging x to the right of y if $x > y$. (This is a non-rigorous arrangement, since we have not yet defined the concept of a line, but this discussion is only intended to motivate the more rigorous propositions below.) Inside the rationals we have the integers, which are thus also arranged on the line. Now we work out how the rationals are arranged with respect to the integers.

Proposition 4.4.1 (Interspersion of integers by rationals). *Let x be a rational number. Then there exists an integer n such that $n \leq x < n+1$. In fact, this integer is unique (i.e., for each x there is only one n for which $n \leq x < n+1$). In particular, there exists a natural number N such that $N > x$ (i.e., there is no such thing as a rational number which is larger than all the natural numbers).*

Remark 4.4.2. The integer n for which $n \leq x < n+1$ is sometimes referred to as the *integer part* of x and is sometimes denoted $n = \lfloor x \rfloor$.

Proof. See Exercise 4.4.1. □

Also, between every two rational numbers there is at least one additional rational:

Proposition 4.4.3 (Interspersion of rationals by rationals). *Given any two rationals x and y such that $x < y$, there exists a third rational z such that $x < z < y$.*

Proof. We set $z := (x + y)/2$. Since $x < y$, and $1/2 = 1/2$ is positive, we have from Proposition 4.2.9 that $x/2 < y/2$. If we add $y/2$ to both sides using Proposition 4.2.9 we obtain $x/2 + y/2 < y/2 + y/2$, i.e., $z < y$. If we instead add $x/2$ to both sides we obtain $x/2 + x/2 < y/2 + x/2$, i.e., $x < z$. Thus $x < z < y$ as desired. \square

Despite the rationals having this denseness property, they are still incomplete; there are still an infinite number of “gaps” or “holes” between the rationals, although this denseness property does ensure that these holes are in some sense infinitely small. For instance, we will now show that the rational numbers do not contain any square root of two.

Proposition 4.4.4. *There does not exist any rational number x for which $x^2 = 2$.*

Proof. We only give a sketch of a proof; the gaps will be filled in Exercise 4.4.3. Suppose for contradiction that we had a rational number x for which $x^2 = 2$. Clearly x is not zero. We may assume that x is positive, for if x were negative then we could just replace x by $-x$ (since $x^2 = (-x)^2$). Thus $x = p/q$ for some positive integers p, q , so $(p/q)^2 = 2$, which we can rearrange as $p^2 = 2q^2$. Define a natural number p to be *even* if $p = 2k$ for some natural number k , and *odd* if $p = 2k + 1$ for some natural number k . Every natural number is either even or odd, but not both (why?). If p is odd, then p^2 is also odd (why?), which contradicts $p^2 = 2q^2$. Thus p is even, i.e., $p = 2k$ for some natural number k . Since p is positive, k must also be positive. Inserting $p = 2k$ into $p^2 = 2q^2$ we obtain $4k^2 = 2q^2$, so that $q^2 = 2k^2$.

To summarize, we started with a pair (p, q) of positive integers such that $p^2 = 2q^2$, and ended up with a pair (q, k) of positive integers such that $q^2 = 2k^2$. Since $p^2 = 2q^2$, we have $q < p$ (why?). If we rewrite $p' := q$ and $q' := k$, we thus can pass from one solution (p, q) to the equation $p^2 = 2q^2$ to a new solution (p', q') to the same equation which has a smaller value of p . But then we can repeat this procedure again and again, obtaining a

sequence (p'', q'') , (p''', q''') , etc. of solutions to $p^2 = 2q^2$, each one with a smaller value of p than the previous, and each one consisting of positive integers. But this contradicts the principle of infinite descent (see Exercise 4.4.2). This contradiction shows that we could not have had a rational x for which $x^2 = 2$. \square

On the other hand, we can get rational numbers which are arbitrarily close to a square root of 2:

Proposition 4.4.5. *For every rational number $\varepsilon > 0$, there exists a non-negative rational number x such that $x^2 < 2 < (x + \varepsilon)^2$.*

Proof. Let $\varepsilon > 0$ be rational. Suppose for contradiction that there is no non-negative rational number x for which $x^2 < 2 < (x + \varepsilon)^2$. This means that whenever x is non-negative and $x^2 < 2$, we must also have $(x + \varepsilon)^2 < 2$ (note that $(x + \varepsilon)^2$ cannot equal 2, by Proposition 4.4.4). Since $0^2 < 2$, we thus have $\varepsilon^2 < 2$, which then implies $(2\varepsilon)^2 < 2$, and indeed a simple induction shows that $(n\varepsilon)^2 < 2$ for every natural number n . (Note that $n\varepsilon$ is non-negative for every natural number n - why?) But, by Proposition 4.4.1 we can find an integer n such that $n > 2/\varepsilon$, which implies that $n\varepsilon > 2$, which implies that $(n\varepsilon)^2 > 4 > 2$, contradicting the claim that $(n\varepsilon)^2 < 2$ for all natural numbers n . This contradiction gives the proof. \square

Example 4.4.6. If³ $\varepsilon = 0.001$, we can take $x = 1.414$, since $x^2 = 1.999396$ and $(x + \varepsilon)^2 = 2.002225$.

Proposition 4.4.5 indicates that, while the set \mathbf{Q} of rationals do not actually have $\sqrt{2}$ as a member, we can get as close as we wish to $\sqrt{2}$. For instance, the sequence of rationals

$$1.4, 1.41, 1.414, 1.4142, 1.41421, \dots$$

seem to get closer and closer to $\sqrt{2}$, as their squares indicate:

$$1.96, 1.9881, 1.99396, 1.99996164, 1.9999899241, \dots$$

³We will use the decimal system for defining terminating decimals, for instance 1.414 is defined to equal the rational number 1414/1000. We defer the formal discussion on the decimal system to an Appendix (§13).

Thus it seems that we can create a square root of 2 by taking a “limit” of a sequence of rationals. This is how we shall construct the real numbers in the next chapter. (There is another way to do so, using something called “Dedekind cuts”, which we will not pursue here. One can also proceed using infinite decimal expansions, but there are some sticky issues when doing so, e.g., one has to make $0.999\dots$ equal to $1.000\dots$, and this approach, despite being the most familiar, is actually *more* complicated than other approaches; see the Appendix §13.)

Exercise 4.4.1. Prove Proposition 4.4.1. (Hint: use Proposition 2.3.9).

Exercise 4.4.2. A definition: a sequence a_0, a_1, a_2, \dots of numbers (natural numbers, integers, rationals, or reals) is said to be in *infinite descent* if we have $a_n > a_{n+1}$ for all natural numbers n (i.e., $a_0 > a_1 > a_2 > \dots$).

- (a) Prove the *principle of infinite descent*: that it is not possible to have a sequence of *natural numbers* which is in infinite descent. (Hint: Assume for contradiction that you can find a sequence of natural numbers which is in infinite descent. Since all the a_n are natural numbers, you know that $a_n \geq 0$ for all n . Now use induction to show in fact that $a_n \geq k$ for all $k \in \mathbf{N}$ and all $n \in \mathbf{N}$, and obtain a contradiction.)
- (b) Does the principle of infinite descent work if the sequence a_1, a_2, a_3, \dots is allowed to take integer values instead of natural number values? What about if it is allowed to take positive rational values instead of natural numbers? Explain.

Exercise 4.4.3. Fill in the gaps marked (why?) in the proof of Proposition 4.4.4.

Chapter 12

Appendix: the basics of mathematical logic

The purpose of this appendix is to give a quick introduction to *mathematical logic*, which is the language one uses to conduct rigorous mathematical proofs. Knowing how mathematical logic works is also very helpful for understanding the mathematical way of thinking, which once mastered allows you to approach mathematical concepts and problems in a clear and confident way - including many of the proof-type questions in this text.

Writing logically is a very useful skill. It is somewhat related to, but not the same as, writing clearly, or efficiently, or convincingly, or informatively; ideally one would want to do all of these at once, but sometimes one has to make compromises (though with practice you'll be able to achieve more of your writing objectives concurrently). Thus a logical argument may sometimes look unwieldy, excessively complicated, or otherwise appear unconvincing. The big advantage of writing logically, however, is that one can be absolutely sure that your conclusion will be correct, as long as all your hypotheses were correct and your steps were logical; using other styles of writing one can be reasonably convinced that something is true, but there is a difference between being convinced and being *sure*.

Being logical is not the only desirable trait in writing, and in fact sometimes it gets in the way; mathematicians for instance often resort to short informal arguments which are not logically rigorous when they want to convince other mathematicians of a

statement without doing through all of the long details, and the same is true of course for non-mathematicians as well. So saying that a statement or argument is “not logical” is not necessarily a bad thing; there are often many situations when one has good reasons to not be emphatic about being logical. However, one should be aware of the distinction between logical reasoning and more informal means of argument, and not try to pass off an illogical argument as being logically rigorous. In particular, if an exercise is asking for a proof, then it is expecting you to be logical in your answer.

Logic is a skill that needs to be learnt like any other, but this skill is also innate to all of you - indeed, you probably use the laws of logic unconsciously in your everyday speech and in your own internal (non-mathematical) reasoning. However, it does take a bit of training and practice to recognize this innate skill and to apply it to abstract situations such as those encountered in mathematical proofs. Because logic is innate, the laws of logic that you learn should *make sense* - if you find yourself having to memorize one of the principles or laws of logic here, without feeling a mental “click” or comprehending why that law should work, then you will probably *not* be able to use that law of logic correctly and effectively in practice. So, *please* don’t study this appendix the way you might cram before a final - that is going to be useless. Instead, **put away your highlighter pen**, and *read* and *understand* this appendix rather than merely *studying* it!

12.1 Mathematical statements

Any mathematical argument proceeds in a sequence of *mathematical statements*. These are precise statements concerning various mathematical objects (numbers, vectors, functions, etc.) and relations between them (addition, equality, differentiation, etc.). These objects can either be constants or variables; more on this later. Statements¹ are either true or false.

¹More precisely, statements with no free variables are either true or false. We shall discuss free variables later on in this appendix.

Example 12.1.1. $2 + 2 = 4$ is a true statement; $2 + 2 = 5$ is a false statement.

Not every combination of mathematical symbols is a statement. For instance,

$$= 2 + +4 = - = 2$$

is not a statement; we sometimes call it *ill-formed* or *ill-defined*. The statements in the previous example are *well-formed* or *well-defined*. Thus well-formed statements can be either true or false; ill-formed statements are considered to be neither true nor false (in fact, they are usually not considered statements at all). A more subtle example of an ill-formed statement is

$$0/0 = 1;$$

division by zero is undefined, and so the above statement is ill-formed. A logical argument should not contain any ill-formed statements, thus for instance if an argument uses a statement such as $x/y = z$, it needs to first ensure that y is not equal to zero. Many purported proofs of “ $0=1$ ” or other false statements rely on overlooking this “statements must be well-formed” criterion.

Many of you have probably written ill-formed or otherwise inaccurate statements in your mathematical work, while intending to mean some other, well-formed and accurate statement. To a certain extent this is permissible - it is similar to misspelling some words in a sentence, or using a slightly inaccurate or ungrammatical word in place of a correct one (“She ran good” instead of “She ran well”). In many cases, the reader (or grader) can detect this mis-step and correct for it. However, it looks unprofessional and suggests that you may not know what you are talking about. And if indeed you actually do not know what you are talking about, and are applying mathematical or logical rules blindly, then writing an ill-formed statement can quickly confuse you into writing more and more nonsense - usually of the sort which receives no credit in grading. So it is important, especially when just learning

a subject, to take care in keeping statements well-formed and precise. Once you have more skill and confidence, of course you can afford once again to speak loosely, because you will know what you are doing and won't be in as much danger of veering off into nonsense.

One of the basic axioms of mathematical logic is that every well-formed statement is either true or false, but not both (though if there are free variables, the truth of a statement may depend on the values of these variables. More on this later). Furthermore, the truth or falsity of a statement is intrinsic to the statement, and does not depend on the opinion of the person viewing the statement (as long as all the definitions and notations are agreed upon, of course). So to prove that a statement is true; it suffices to show that it is not false; to show that a statement is false, it suffices to show that it is not true; this is the principle underlying the powerful technique of *proof by contradiction*. This axiom is viable as long as one is working with precise concepts, for which the truth or falsity can be determined (at least in principle) in an objective and consistent manner. However, if one is working in very non-mathematical situations, then this axiom becomes much more dubious, and so it can be a mistake to apply mathematical logic to non-mathematical situations. (For instance, a statement such as “this rock weighs 52 pounds” is reasonably precise and objective, and so it is fairly safe to use mathematical reasoning to manipulate it, whereas statements such as “this rock is heavy”, “this piece of music is beautiful” or “God exists” are much more problematic. So while mathematical logic is a very useful and powerful tool, it still does have some limitations of applicability.) One can still attempt to apply logic (or principles similar to logic) in these cases (for instance, by creating a *mathematical model* of a real-life phenomenon), but this is now science or philosophy, not mathematics, and we will not discuss it further here.

Remark 12.1.2. There are other models of logic which attempts to deal with statements that are not definitely true or definitely false, such as modal logics, intuitionist logics, or fuzzy logics, but these are definitely in the realm of logic and philosophy and thus

well beyond the scope of this text.

Being true is different from being *useful* or *efficient*. For instance, the statement

$$2 = 2$$

is true but unlikely to be very useful. The statement

$$4 \leq 4$$

is also true, but not very efficient (the statement $4 = 4$ is more precise). It may also be that a statement may be false yet still be useful, for instance

$$\pi = 22/7$$

is false, but is still useful as a first approximation. In mathematical reasoning, we only concern ourselves with truth rather than usefulness or efficiency; the reason is that truth is objective (everybody can agree on it) and we can deduce true statements from precise rules, whereas usefulness and efficiency are to some extent matters of opinion, and do not follow precise rules. Also, even if some of the individual steps in an argument may not seem very useful or efficient, it is still possible (indeed, quite common) for the final conclusion to be quite non-trivial (i.e., not obviously true) and useful.

Statements are different from *expressions*. Statements are true or false; expressions are a sequence of mathematical sequence which produces some mathematical object (a number, matrix, function, set, etc.) as its value. For instance

$$2 + 3 * 5$$

is an expression, not a statement; it produces a number as its value. Meanwhile,

$$2 + 3 * 5 = 17$$

is a statement, not an expression. Thus it does not make any sense to ask whether $2 + 3 * 5$ is true or false. As with statements, expressions can be well-defined or ill-defined; $2 + 3/0$, for instance,

is ill-defined. More subtle examples of ill-defined expressions arise when, for instance, attempting to add a vector to a matrix, or evaluating a function outside of its domain, e.g., $\sin^{-1}(2)$.

One can make statements out of expressions by using *relations* such as $=$, $<$, \geq , \in , \subset , etc. or by using *properties* (such as “is prime”, “is continuous”, “is invertible”, etc.) For instance, “ $30 + 5$ is prime” is a statement, as is “ $30 + 5 \leq 42 - 7$ ”. Note that mathematical statements are allowed to contain English words.

One can make a *compound statement* from more primitive statements by using *logical connectives* such as and, or, not, if-then, if-and-only-if, and so forth. We give some examples below, in decreasing order of intuitiveness.

Conjunction. If X is a statement and Y is a statement, the statement “ X and Y ” is true if X and Y are both true, and is false otherwise. For instance, “ $2 + 2 = 4$ and $3 + 3 = 6$ ” is true, while “ $2 + 2 = 4$ and $3 + 3 = 5$ ” is not. Another example: “ $2 + 2 = 4$ and $2 + 2 = 4$ ” is true, even if it is a bit redundant; logic is concerned with truth, not efficiency.

Due to the expressiveness of the English language, one can reword the statement “ X and Y ” in many ways, e.g., “ X and also Y ”, or “Both X and Y are true”, etc. Interestingly, the statement “ X , but Y ” is logically the same statement as “ X and Y ”, but they have different connotations (both statements affirm that X and Y are both true, but the first version suggests that X and Y are in contrast to each other, while the second version suggests that X and Y support each other). Again, logic is about truth, not about connotations or suggestions.

Disjunction. If X is a statement and Y is a statement, the statement “ X or Y ” is true if either X or Y is true, or both. For instance, “ $2 + 2 = 4$ or $3 + 3 = 5$ ” is true, but “ $2 + 2 = 5$ or $3 + 3 = 5$ ” is not. Also “ $2 + 2 = 4$ or $3 + 3 = 6$ ” is true (even if it is a bit inefficient; it would be a stronger statement to say “ $2 + 2 = 4$ and $3 + 3 = 6$ ”). Thus by default, the word “or” in mathematical logic defaults to *inclusive or*. The reason we do this is that with inclusive or, to verify “ X or Y ”, it suffices to verify that just one of X or Y is true; we don’t need to show that the

other one is false. So we know, for instance, that “ $2 + 2 = 4$ or $2353 + 5931 = 7284$ ” is true without having to look at the second equation. As in the previous discussion, the statement “ $2 + 2 = 4$ or $2 + 2 = 4$ ” is true, even if it is highly inefficient.

If one really does want to use exclusive or, use a statement such as “Either X or Y is true, but not both” or “Exactly one of X or Y is true”. Exclusive or does come up in mathematics, but nowhere near as often as inclusive or.

Negation. The statement “ X is not true” or “ X is false”, or “It is not the case that X ”, is called the *negation* of X , and is true if and only if X is false, and is false if and only if X is true. For instance, the statement “It is not the case that $2 + 2 = 5$ ” is a true statement. Of course we could abbreviate this statement to “ $2 + 2 \neq 5$ ”.

Negations convert “and” into “or”. For instance, the negation of “Jane Doe has black hair and Jane Doe has blue eyes” is “Jane Doe doesn’t have black hair or doesn’t have blue eyes”, *not* “Jane Doe doesn’t have black hair and doesn’t have blue eyes” (can you see why?). Similarly, if x is an integer, the negation of “ x is even and non-negative” is “ x is odd or negative”, not “ x is odd and negative” (Note how it is important here that or is inclusive rather than exclusive). Or the negation of “ $x \geq 2$ and $x \leq 6$ ” (i.e., “ $2 \leq x \leq 6$ ”) is “ $x < 2$ or $x > 6$ ”, not “ $x < 2$ and $x > 6$ ” or “ $2 < x > 6$ ”.

Similarly, negations convert “or” into “and”. The negation of “John Doe has brown hair or black hair” is “John Doe does not have brown hair and does not have black hair”, or equivalently “John Doe has neither brown nor black hair”. If x is a real number, the negation of “ $x \geq 1$ or $x \leq -1$ ” is “ $x < 1$ and $x > -1$ ” (i.e., $-1 < x < 1$).

It is quite possible that a negation of a statement will produce a statement which could not possibly be true. For instance, if x is an integer, the negation of “ x is either even or odd” is “ x is neither even nor odd”, which cannot possibly be true. Remember, though, that even if a statement is false, it is still a statement, and it is definitely possible to arrive at a true statement using an argument

which at times involves false statements. (Proofs by contradiction, for instance, fall into this category. Another example is proof by dividing into cases. If one divides into three mutually exclusive cases, Case 1, Case 2, and Case 3, then at any given time two of the cases will be false and only one will be true, however this does not necessarily mean that the proof as a whole is incorrect or that the conclusion is false.)

Negations are sometimes unintuitive to work with, especially if there are multiple negations; a statement such as “It is not the case that either x is not odd, or x is not larger than or equal to 3, but not both” is not particularly pleasant to use. Fortunately, one rarely has to work with more than one or two negations at a time, since often negations cancel each other. For instance, the negation of “ X is not true” is just “ X is true”, or more succinctly just “ X ”. Of course one should be careful when negating more complicated expressions because of the switching of “and” and “or”, and similar issues.

If and only if (iff). If X is a statement, and Y is a statement, we say that “ X is true if and only if Y is true”, whenever X is true, Y has to be also, and whenever Y is true, X has to be also (i.e., X and Y are “equally true”). Other ways of saying the same thing are “ X and Y are logically equivalent statements”, or “ X is true iff Y is true”, or “ $X \leftrightarrow Y$ ”. Thus for instance, if x is a real number, then the statement “ $x = 3$ if and only if $2x = 6$ ” is true: this means that whenever $x = 3$ is true, then $2x = 6$ is true, and whenever $2x = 6$ is true, then $x = 3$ is true. On the other hand, the statement “ $x = 3$ if and only if $x^2 = 9$ ” is false; while it is true that whenever $x = 3$ is true, $x^2 = 9$ is also true, it is not the case that whenever $x^2 = 9$ is true, that $x = 3$ is also automatically true (think of what happens when $x = -3$).

Statements that are equally true, are also equally false: if X and Y are logically equivalent, and X is false, then Y has to be false also (because if Y were true, then X would also have to be true). Conversely, any two statements which are equally false will also be logically equivalent. Thus for instance $2 + 2 = 5$ if and only if $4 + 4 = 10$.

Sometimes it is of interest to show that more than two statements are logically equivalent; for instance, one might want to assert that three statements X , Y , and Z are all logically equivalent. This means whenever one of the statements is true, then all of the statements are true; and it also means that if one of the statements is false, then all of the statements are false. This may seem like a lot of logical implications to prove, but in practice, once one demonstrates enough logical implications between X , Y , and Z , one can often conclude all the others and conclude that they are all logically equivalent. See for instance Exercises 12.1.5, 12.1.6.

Exercise 12.1.1. What is the negation of the statement “either X is true, or Y is true, but not both”?

Exercise 12.1.2. What is the negation of the statement “ X is true if and only if Y is true”? (There may be multiple ways to phrase this negation).

Exercise 12.1.3. Suppose that you have shown that whenever X is true, then Y is true, and whenever X is false, then Y is false. Have you now demonstrated that X and Y are logically equivalent? Explain.

Exercise 12.1.4. Suppose that you have shown that whenever X is true, then Y is true, and whenever Y is false, then X is false. Have you now demonstrated that X is true if and only if Y is true? Explain.

Exercise 12.1.5. Suppose you know that X is true if and only if Y is true, and you know that Y is true if and only if Z is true. Is this enough to show that X, Y, Z are all logically equivalent? Explain.

Exercise 12.1.6. Suppose you know that whenever X is true, then Y is true; that whenever Y is true, then Z is true; and whenever Z is true, then X is true. Is this enough to show that X, Y, Z are all logically equivalent? Explain.

12.2 Implication

Now we come to the least intuitive of the commonly used logical connectives - implication. If X is a statement, and Y is a statement, then “if X , then Y ” is the implication from X to Y ; it is also written “when X is true, Y is true”, or “ X implies Y ” or “ Y is true when X is” or “ X is true only if Y is true” (this last one takes a bit of mental effort to see). What this statement “if X , then Y ” means depends on whether X is true or false. If X is true, then “if X , then Y ” is true when Y is true, and false when Y is false. If however X is false, then “if X , then Y ” is *always* true, regardless of whether Y is true or false! To put it another way, when X is true, the statement “if X , then Y ” implies that Y is true. But when X is false, the statement “if X , then Y ” offers no information about whether Y is true or not; the statement is true, but *vacuous* (i.e., does not convey any new information beyond the fact that the hypothesis is false).

Examples 12.2.1. If x is an integer, then the statement “If $x = 2$, then $x^2 = 4$ ” is true, regardless of whether x is actually equal to 2 or not (though this statement is only likely to be useful when x is equal to 2). This statement does not assert that x is equal to 2, and does not assert that x^2 is equal to 4, but it does assert that when and if x is equal to 2, then x^2 is equal to 4. If x is not equal to 2, the statement is still true but offers no conclusion on x or x^2 .

Some special cases of the above implication: the implication “If $2 = 2$, then $2^2 = 4$ ” is true (true implies true). The implication “If $3 = 2$, then $3^2 = 4$ ” is true (false implies false). The implication “If $-2 = 2$, then $(-2)^2 = 4$ ” is true (false implies true). The latter two implications are considered vacuous - they do not offer any new information since their hypothesis is false. (Nevertheless, it is still possible to employ vacuous implications to good effect in a proof - a vacuously true statement is still true. We shall see one such example shortly).

As we see, the falsity of the hypothesis does not destroy the truth of an implication, in fact it is just the opposite! (When

a hypothesis is false, the implication is automatically true.) The only way to disprove an implication is to show that the hypothesis is true while the conclusion is false. Thus “If $2 + 2 = 4$, then $4 + 4 = 2$ ” is a false implication. (True does not imply false.)

One can also think of the statement “if X , then Y ” as “ Y is *at least as true as* X ” - if X is true, then Y also has to be true, but if X is false, Y could be as false as X , but it could also be true. This should be compared with “ X if and only if Y ”, which asserts that X and Y are *equally true*.

Vacuously true implications are often used in ordinary speech, sometimes without knowing that the implication is vacuous. A somewhat frivolous example is “If wishes were wings, then pigs would fly”. (The statement “hell freezes over” is also a popular choice for a false hypothesis.) A more serious one is “If John had left work at 5pm, then he would be here by now.” This kind of statement is often used in a situation in which the conclusion and hypothesis are both false; but the implication is still true regardless. This statement, by the way, can be used to illustrate the technique of proof by contradiction: if you believe that “If John had left work at 5pm, then he would be here by now”, and you also know that “John is not here by now”, then you can conclude that “John did not leave work at 5pm”, because John leaving work at 5pm would lead to a contradiction. Note how a vacuous implication can be used to derive a useful truth.

To summarize, implications are sometimes vacuous, but this is not actually a problem in logic, since these implications are still true, and vacuous implications can still be useful in logical arguments. In particular, one can safely use statements like “If X , then Y ” without necessarily having to worry about whether the hypothesis X is actually true or not (i.e., whether the implication is vacuous or not).

Implications can also be true even when there is no causal link between the hypothesis and conclusion. The statement “If $1 + 1 = 2$, then Washington D.C. is the capital of the United States” is true (true implies true), although rather odd; the statement “If $2 + 2 = 3$, then New York is the capital of the United

States” is similarly true (false implies false). Of course, such a statement may be unstable (the capital of the United States may one day change, while $1 + 1$ will always remain equal to 2) but it is true, at least for the moment. While it is possible to use acausal implications in a logical argument, it is not recommended as it can cause unneeded confusion. (Thus, for instance, while it is true that a false statement can be used to imply any other statement, true or false, doing so arbitrarily would probably not be helpful to the reader.)

To prove an implication “If X , then Y ”, the usual way to do this is to first assume that X is true, and use this (together with whatever other facts and hypotheses you have) to deduce Y . This is still a valid procedure even if X later turns out to be false; the implication does not guarantee anything about the truth of X , and only guarantees the truth of Y conditionally on X first being true. For instance, the following is a valid proof of a true Proposition, even though both hypothesis and conclusion of the Proposition are false:

Proposition 12.2.2. *If $2 + 2 = 5$, then $4 = 10 - 4$.*

Proof. Assume $2 + 2 = 5$. Multiplying both sides by 2, we obtain $4 + 4 = 10$. Subtracting 4 from both sides, we obtain $4 = 10 - 4$ as desired. \square

On the other hand, a common error is to prove an implication by first assuming the *conclusion* and then arriving at the hypothesis. For instance, the following Proposition is correct, but the proof is not:

Proposition 12.2.3. *Suppose that $2x + 3 = 7$. Show that $x = 2$.*

Proof. (Incorrect) $x = 2$; so $2x = 4$; so $2x + 3 = 7$. \square

When doing proofs, it is important that you are able to distinguish the hypothesis from the conclusion; there is a danger of getting hopelessly confused if this distinction is not clear.

Here is a short proof which uses implications which are possibly vacuous.

Theorem 12.2.4. *Suppose that n is an integer. Then $n(n + 1)$ is an even integer.*

Proof. Since n is an integer, n is even or odd. If n is even, then $n(n + 1)$ is also even, since any multiple of an even number is even. If n is odd, then $n + 1$ is even, which again implies that $n(n + 1)$ is even. Thus in either case $n(n + 1)$ is even, and we are done. \square

Note that this proof relied on two implications: “if n is even, then $n(n + 1)$ is even”, and “if n is odd, then $n(n + 1)$ is even”. Since n cannot be both odd and even, at least one of these implications has a false hypothesis and is therefore vacuous. Nevertheless, both these implications are true, and one needs *both* of them in order to prove the theorem, because we don’t know in advance whether n is even or odd. And even if we did, it might not be worth the trouble to check it. For instance, as a special case of this Theorem we immediately know

Corollary 12.2.5. *Let $n = (253 + 142) * 123 - (423 + 198)^{342} + 538 - 213$. Then $n(n + 1)$ is an even integer.*

In this particular case, one can work out exactly which parity n is - even or odd - and then use only one of the two implications in above the Theorem, discarding the vacuous one. This may seem like it is more efficient, but it is a false economy, because one then has to determine what parity n is, and this requires a bit of effort - more effort than it would take if we had just left both implications, including the vacuous one, in the argument. So, somewhat paradoxically, the inclusion of vacuous, false, or otherwise “useless” statements in an argument can actually *save* you effort in the long run! (I’m not suggesting, of course, that you ought to pack your proofs with lots of time-wasting and irrelevant statements; all I’m saying here is that you need not be unduly concerned that some hypotheses in your argument might not be correct, as long as your argument is still structured to give the correct conclusion regardless of whether those hypotheses were true or false.)

The statement “If X , then Y ” is not the same as “If Y , then X ”; for instance, while “If $x = 2$, then $x^2 = 4$ ” is true, “If $x^2 = 4$, then $x = 2$ ” can be false if x is equal to -2 . These two statements are called *converses* of each other; thus the converse of a true implication is not necessarily another true implication. We use the statement “ X if and only if Y ” to denote the statement that “If X , then Y ; and if Y , then X ”. Thus for instance, we can say that $x = 2$ if and only if $2x = 4$, because if $x = 2$ then $2x = 4$, while if $2x = 4$ then $x = 2$. One way of thinking about an if-and-only-if statement is to view “ X if and only if Y ” as saying that X is just as true as Y ; if one is true then so is the other, and if one is false, then so is the other. For instance, the statement “If $3 = 2$, then $6 = 4$ ” is true, since both hypothesis and conclusion are false. (Under this view, “If X , then Y ” can be viewed as a statement that Y is at least as true as X .) Thus one could say “ X and Y are equally true” instead of “ X if and only if Y ”.

Similarly, the statement “If X is true, then Y is true” is NOT the same as “If X is false, then Y is false”. Saying that “if $x = 2$, then $x^2 = 4$ ” does not imply that “if $x \neq 2$, then $x^2 \neq 4$ ”, and indeed we have $x = -2$ as a counterexample in this case. If-then statements are not the same as if-and-only-if statements. (If we knew that “ X is true if and only if Y is true”, then we would also know that “ X is false if and only if Y is false”.) The statement “If X is false, then Y is false” is sometimes called the *inverse* of “If X is true, then Y is true”; thus the inverse of a true implication is not necessarily a true implication.

If you know that “If X is true, then Y is true”, then it is also true that “If Y is false, then X is false” (because if Y is false, then X can’t be true, since that would imply Y is true, a contradiction.) For instance, if we knew that “If $x = 2$, then $x^2 = 4$ ”, then we also know that “If $x^2 \neq 4$, then $x \neq 2$ ”. Or if we knew “If John had left work at 5pm, he would be here by now”, then we also know “If John isn’t here now, then he could not have left work at 5pm”. The statement “If Y is false, then X is false” is known as the *contrapositive* of “If X , then Y ” and both statements are equally true.

In particular, if you know that X implies something which is known to be false, then X itself must be false. This is the idea behind *proof by contradiction* or *reductio ad absurdum*: to show something must be false, assume first that it is true, and show that this implies something which you know to be false (e.g., that a statement is simultaneously true and not true). For instance:

Proposition 12.2.6. *Let x be a positive number such that $\sin(x) = 1$. Then $x \geq \pi/2$.*

Proof. Suppose for contradiction that $x < \pi/2$. Since x is positive, we thus have $0 < x < \pi/2$. Since $\sin(x)$ is increasing for $0 < x < \pi/2$, and $\sin(0) = 0$ and $\sin(\pi/2) = 1$, we thus have $0 < \sin(x) < 1$. But this contradicts the hypothesis that $\sin(x) = 1$. Hence $x \geq \pi/2$. \square

Note that one feature of proof by contradiction is that at some point in the proof you assume a hypothesis (in this case, that $x < \pi/2$) which later turns out to be false. Note however that this does not alter the fact that the argument remains valid, and that the conclusion is true; this is because the ultimate conclusion does not rely on that hypothesis being true (indeed, it relies instead on it being false!).

Proof by contradiction is particularly useful for showing “negative” statements - that X is false, that a is not equal to b , that kind of thing. But the line between positive and negative statements is sort of blurry (is the statement $x \geq 2$ a positive or negative statement? What about its negation, that $x < 2$?) so this is not a hard and fast rule.

Logicians often use special symbols to denote logical connectives; for instance “ X implies Y ” can be written “ $X \implies Y$ ”, “ X is not true” can be written “ $\sim X$ ”, “! X ”, or “ $\neg X$ ”, “ X and Y ” can be written “ $X \wedge Y$ ” or “ $X\&Y$ ”, and so forth. But for general-purpose mathematics, these symbols are not often used; English words are often more readable, and don’t take up much more space. Also, using these symbols tends to blur the line between expressions and statements; it’s not as easy to parse

“ $((x = 3) \wedge (y = 5)) \implies (x + y = 8)$ ” as “If $x = 3$ and $y = 5$, then $x + y = 8$ ”. So in general I would not recommend using these symbols (except possibly for \implies , which is a very intuitive symbol).

12.3 The structure of proofs

To prove a statement, one often starts by assuming the hypothesis and working one’s way toward a conclusion; this is the *direct* approach to proving a statement. Such a proof might look something like the following:

Proposition 12.3.1. *A implies B.*

Proof. Assume A is true. Since A is true, C is true. Since C is true, D is true. Since D is true, B is true, as desired. \square

An example of such a direct approach is

Proposition 12.3.2. *If $x = \pi$, then $\sin(x/2) + 1 = 2$.*

Proof. Let $x = \pi$. Since $x = \pi$, we have $x/2 = \pi/2$. Since $x/2 = \pi/2$, we have $\sin(x/2) = 1$. Since $\sin(x/2) = 1$, we have $\sin(x/2) + 1 = 2$. \square

Note what we did here was started at the hypothesis and moved steadily from there toward a conclusion. It is also possible to work backwards from the conclusion and seeing what it would take to imply it. For instance, a typical proof of Proposition 12.3.1 of this sort might look like the following:

Proof. To show B , it would suffice to show D . Since C implies D , we just need to show C . But C follows from A . \square

As an example of this, we give another proof of Proposition 12.3.2:

Proof. To show $\sin(x/2) + 1 = 2$, it would suffice to show that $\sin(x/2) = 1$. Since $x/2 = \pi/2$ would imply $\sin(x/2) = 1$, we just need to show that $x/2 = \pi/2$. But this follows since $x = \pi$. \square

Logically speaking, the above two proofs of Proposition 12.3.2 are the same, just arranged differently. Note how this proof style is different from the (incorrect) approach of starting with the conclusion and seeing what it would imply (as in Proposition 12.2.3); instead, we start with the conclusion and see what would imply it.

Another example of a proof written in this backwards style is the following:

Proposition 12.3.3. *Let $0 < r < 1$ be a real number. Then the series $\sum_{n=1}^{\infty} nr^n$ is convergent.*

Proof. To show this series is convergent, it suffices by the ratio test to show that the ratio

$$\left| \frac{r^{n+1}(n+1)}{r^n n} \right| = r \frac{n+1}{n}$$

converges to something less than 1 as $n \rightarrow \infty$. Since r is already less than 1, it will be enough to show that $\frac{n+1}{n}$ converges to 1. But since $\frac{n+1}{n} = 1 + \frac{1}{n}$, it suffices to show that $\frac{1}{n} \rightarrow 0$. But this is clear since $n \rightarrow \infty$. \square

One could also do any combination of moving forwards from the hypothesis and backwards from the conclusion. For instance, the following would be a valid proof of Proposition 12.3.1:

Proof. To show B , it would suffice to show D . So now let us show D . Since we have A by hypothesis, we have C . Since C implies D , we thus have D as desired. \square

Again, from a logical point of view this is exactly the same proof as before. Thus there are many ways to write the same proof down; how you do so is up to you, but certain ways of writing proofs are more readable and natural than others, and different arrangements tend to emphasize different parts of the argument. (Of course, when you are just starting out doing mathematical proofs, you're generally happy to get *some* proof of a result, and don't care so much about getting the "best" arrangement of that

proof; but the point here is that a proof can take many different forms.)

The above proofs were pretty simple because there was just one hypothesis and one conclusion. When there are multiple hypotheses and conclusions, and the proof splits into cases, then proofs can get more complicated. For instance a proof might look as tortuous as this:

Proposition 12.3.4. *Suppose that A and B are true. Then C and D are true.*

Proof. Since A is true, E is true. From E and B we know that F is true. Also, in light of A , to show D it suffices to show G . There are now two cases: H and I . If H is true, then from F and H we obtain C , and from A and H we obtain G . If instead I is true, then from I we have G , and from I and G we obtain C . Thus in both cases we obtain both C and G , and hence C and D . \square

Incidentally, the above proof could be rearranged into a much tidier manner, but you at least get the idea of how complicated a proof could become. To show an implication there are several ways to proceed: you can work forward from the hypothesis; you can work backward from the conclusion; or you can divide into cases in the hope to split the problem into several easier sub-problems. Another is to argue by contradiction, for instance you can have an argument of the form

Proposition 12.3.5. *Suppose that A is true. Then B is false.*

Proof. Suppose for contradiction that B is true. This would imply that C is true. But since A is true, this implies that D is true; which contradicts C . Thus B must be false. \square

As you can see, there are several things to try when attempting a proof. With experience, it will become clearer which approaches are likely to work easily, which ones will probably work but require much effort, and which ones are probably going to fail. In many cases there is really only one obvious way to proceed. Of course,

there may definitely be multiple ways to approach a problem, so if you see more than one way to begin a problem, you can just try whichever one looks the easiest, but be prepared to switch to another approach if it begins to look hopeless.

Also, it helps when doing a proof to keep track of which statements are *known* (either as hypotheses, or deduced from the hypotheses, or coming from other theorems and results), and which statements are *desired* (either the conclusion, or something which would imply the conclusion, or some intermediate claim or lemma which will be useful in eventually obtaining the conclusion). Mixing the two up is almost always a bad idea, and can lead to one getting hopelessly lost in a proof.

12.4 Variables and quantifiers

One can get quite far in logic just by starting with primitive statements (such as “ $2+2 = 4$ ” or “John has black hair”), then forming compound statements using logical connectives, and then using various laws of logic to pass from one’s hypotheses to one’s conclusions; this is known as *propositional logic* or *Boolean logic*. (It is possible to list a dozen or so such laws of propositional logic, which are sufficient to do everything one wants to do, but I have deliberately chosen not to do so here, because you might then be tempted to memorize that list, and that is NOT how one should learn how to do logic, unless one happens to be a computer or some other non-thinking device. However, if you really are curious as to what the formal laws of logic are, look up “laws of propositional logic” or something similar in the library or on the internet.)

However, to do mathematics, this level of logic is insufficient, because it does not incorporate the fundamental concept of *variables* - those familiar symbols such as x or n which denote various quantities which are unknown, or set to some value, or assumed to obey some property. Indeed we have already sneaked in some of these variables in order to illustrate some of the concepts in propositional logic (mainly because it gets boring after a while to

talk endlessly about variable-free statements such as $2 + 2 = 4$ or “Jane has black hair”). *Mathematical logic* is thus the same as propositional logic but with the additional ingredient of variables added.

A *variable* is a symbol, such as n or x , which denotes a certain type of mathematical object - an integer, a vector, a matrix, that kind of thing. In almost all circumstances, the type of object that the variable is should be declared, otherwise it will be difficult to make well-formed statements using it. (There are a very few number of true statements one can make about variables which can be of absolutely any type. For instance, given a variable x of any type whatsoever, it is true that $x = x$, and if we also know that $x = y$, then we can conclude that $y = x$. But one cannot say, for instance, that $x + y = y + x$, until we know what type of objects x and y are and whether they support the operation of addition; for instance, the above statement is ill-formed if x is a matrix and y is a vector. Thus if one actually wants to do some useful mathematics, then every variable should have an explicit type.)

One can form expressions and statements involving variables, for instance, if x is a real variable (i.e., a variable which is a real number), $x + 3$ is an expression, and $x + 3 = 5$ is a statement. But now the truth of a statement may depend on the value of the variables involved; for instance the statement $x + 3 = 5$ is true if x is equal to 2, but is false if x is not equal to 2. Thus the truth of a statement involving a variable may depend on the *context* of the statement - in this case, it depends on what x is supposed to be. (This is a modification of the rule for propositional logic, in which all statements have a definite truth value.)

Sometimes we do not set a variable to be anything (other than specifying its type). Thus, we could consider the statement $x + 3 = 5$ where x is an unspecified real number. In such a case we call this variable a *free variable*; thus we are considering $x + 3 = 5$ with x a free variable. Statements with free variables might not have a definite truth value, as they depend on an unspecified variable. For instance, we have already remarked that $x + 3 = 5$ does not

have a definite truth value if x is a free real variable, though of course for each given value of x the statement is either true or false. On the other hand, the statement $(x + 1)^2 = x^2 + 2x + 1$ is true for every real number x , and so we can regard this as a true statement even when x is a free variable.

At other times, we *set* a variable to equal a fixed value, by using a statement such as “Let $x = 2$ ” or “Set x equal to 2”. In that case, the variable is known as a *bound variable*, and statements involving only bound variables and no free variables do have a definite truth value. For instance, if we set $x = 342$, then the statement “ $x + 135 = 477$ ” now has a definite truth value, whereas if x is a free real variable then “ $x + 135 = 477$ ” could be either true or false, depending on what x is. Thus, as we have said before, the truth of a statement such as “ $x + 135 = 477$ ” depends on the context - whether x is free or bound, and if it is bound, what it is bound to.

One can also turn a free variable into a bound variable by using the quantifiers “for all” or “for some”. For instance, the statement

$$(x + 1)^2 = x^2 + 2x + 1$$

is a statement with one free variable x , and need not have a definite truth value, but the statement

$$(x + 1)^2 = x^2 + 2x + 1 \text{ for all real numbers } x$$

is a statement with one bound variable x , and now has a definite truth value (in this case, the statement is true). Similarly, the statement

$$x + 3 = 5$$

has one free variable, and does not have a definite truth value, but the statement

$$x + 3 = 5 \text{ for some real number } x$$

is true, since it is true for $x = 2$. On the other hand, the statement

$$x + 3 = 5 \text{ for all real numbers } x$$

is false, because there are some (indeed, there are many) real numbers x for which $x + 3$ is not equal to 5.

Universal quantifiers. Let $P(x)$ be some statement depending on a free variable x . The statement “ $P(x)$ is true for all x of type T ” means that given any x of type T , the statement $P(x)$ is true regardless of what the exact value of x is. In other words, the statement is the same as saying “if x is of type T , then $P(x)$ is true”. Thus the usual way to prove such a statement is to let x be a free variable of type T (by saying something like “Let x be any object of type T ”), and then proving $P(x)$ for that object. The statement becomes false if one can produce even a single counterexample, i.e., an element x which lies in T but for which $P(x)$ is false. For instance, the statement “ x^2 is greater than x for all positive x ” can be shown to be false by producing a single example, such as $x = 1$ or $x = 1/2$, where x^2 is not greater than x .

On the other hand, producing a single example where $P(x)$ is true will not show that $P(x)$ is true for *all* x . For instance, just because the equation $x + 3 = 5$ has a solution when $x = 2$ does not imply that $x + 3 = 5$ for all real numbers x ; it only shows that $x + 3 = 5$ is true for some real number x . (This is the source of the often-quoted, though somewhat inaccurate, slogan “One cannot prove a statement just by giving an example”. The more precise statement is that one cannot prove a “for all” statement by examples, though one can certainly prove “for some” statements this way, and one can also *disprove* “for all” statements by a single counterexample.)

It occasionally happens that there are in fact no variables x of type T . In that case the statement “ $P(x)$ is true for all x of type T ” is *vacuously true* - it is true but has no content, similar to a vacuous implication. For instance, the statement

$$6 < 2x < 4 \text{ for all } 3 < x < 2$$

is true, and easily proven, but is vacuous. (Such a vacuously true statement can still be useful in an argument, this doesn’t happen very often.)

One can use phrases such as “For every” or “For each” instead of “For all”, e.g., one can rephrase “ $(x + 1)^2 = x^2 + 2x + 1$ for all real numbers x ” as “For each real number x , $(x + 1)^2$ is equal to $x^2 + 2x + 1$ ”. For the purposes of logic these rephrasings are equivalent. The symbol \forall can be used instead of “For all”, thus for instance “ $\forall x \in X : P(x)$ is true” or “ $P(x)$ is true $\forall x \in X$ ” is synonymous with “ $P(x)$ is true for all $x \in X$ ”.

Existential quantifiers The statement “ $P(x)$ is true for some x of type T ” means that there exists at least one x of type T for which $P(x)$ is true, although it may be that there is more than one such x . (One would use a quantifier such as “for exactly one x ” instead of “for some x ” if one wanted both existence and uniqueness of such an x). To prove such a statement it suffices to provide a single example of such an x . For instance, to show that

$$x^2 + 2x - 8 = 0 \text{ for some real number } x$$

all one needs to do is find a single real number x for which $x^2 + 2x - 8 = 0$, for instance $x = 2$ will do. (One could also use $x = -4$, but one doesn’t need to use both.) Note that one has the freedom to select x to be anything you want when proving a for-some statement; this is in contrast to proving a for-all statement, where you have to let x be arbitrary. (One can compare the two statements by thinking of two games between you and an opponent. In the first game, the opponent gets to pick what x is, and then you have to prove $P(x)$; if you can always win this game, then you have proven that $P(x)$ is true for *all* x . In the second game, *you* get to choose what x is, and then you prove $P(x)$; if you can win this game, you have proven that $P(x)$ is true for *some* x .)

Usually, saying something is true for *all* x is much stronger than just saying it is true for *some* x . There is one exception though, if the condition on x is impossible to satisfy, then the for-all statement is vacuously true, but the for-some statement is false. For instance

$$6 < 2x < 4 \text{ for all } 3 < x < 2$$

is true, but

$$6 < 2x < 4 \text{ for some } 3 < x < 2$$

is false.

One can use phrases such as “For at least one” or “There exists ... such that” instead of “For some”. For instance, one can rephrase “ $x^2 + 2x - 8$ for some real number x ” as “There exists a real number x such that $x^2 + 2x - 8$ ”. The symbol \exists can be used instead of “There exists ... such that”, thus for instance “ $\exists x \in X : P(x)$ is true” is synonymous with “ $P(x)$ is true for some $x \in X$ ”.

12.5 Nested quantifiers

One can nest two or more quantifiers together. For instance, consider the statement

For every positive number x , there exists a positive number y such that $y^2 = x$.

What does this statement mean? It means that for each positive number x , the statement

There exists a positive number y such that $y^2 = x$

is true. In other words, one can find a positive square root of x for each positive number x . So the statement is saying that every positive number has a positive square root.

To continue the gaming metaphor, suppose you play a game where your opponent first picks a positive number x , and then you pick a positive number y . You win the game if $y^2 = x$. If you can always win the game regardless of what your opponent does, then you have proven that for every positive x , there exists a positive y such that $y^2 = x$.

Negating a universal statement produces an existential statement. The negation of “All swans are white” is not “All swans are not white”, but rather “There is some swan which is not white”. Similarly, the negation of “For every $0 < x < \pi/2$, we have $\cos(x) \geq 0$ ” is “For some $0 < x < \pi/2$, we have $\cos(x) < 0$ ”, NOT “For every $0 < x < \pi/2$, we have $\cos(x) < 0$ ”.

Negating an existential statement produces a universal statement. The negation of “There exists a black swan” is not “There

exists a swan which is not black”, but rather “All swans are not black”. Similarly, the negation of “There exists a real number x such that $x^2 + x + 1 = 0$ ” is “For every real number x , $x^2 + x + 1 \neq 0$ ”, NOT “There exists a real number x such that $x^2 + x + 1 \neq 0$ ”. (The situation here is very similar to how “and” and “or” behave with respect to negations.)

If you know that a statement $P(x)$ is true for all x , then you can set x to be anything you want, and $P(x)$ will be true for that value of x ; this is what “for all” means. Thus for instance if you know that

$$(x + 1)^2 = x^2 + 2x + 1 \text{ for all real numbers } x,$$

then you can conclude for instance that

$$(\pi + 1)^2 = \pi^2 + 2\pi + 1,$$

or for instance that

$$(\cos(y) + 1)^2 = \cos(y)^2 + 2\cos(y) + 1 \text{ for all real numbers } y$$

(because if y is real, then $\cos(y)$ is also real), and so forth. Thus universal statements are very versatile in their applicability - you can get $P(x)$ to hold for whatever x you wish. Existential statements, by contrast, are more limited; if you know that

$$x^2 + 2x - 8 = 0 \text{ for some real number } x$$

then you cannot simply substitute in any real number you wish, e.g., π , and conclude that $\pi^2 + 2\pi - 8 = 0$. However, you can of course still conclude that $x^2 + 2x - 8 = 0$ for some real number x , it’s just that you don’t get to pick which x it is. (To continue the gaming metaphor, you can make $P(x)$ hold, but your opponent gets to pick x for you, you don’t get to choose for yourself.)

Remark 12.5.1. In the history of logic, quantifiers were formally studied thousands of years before Boolean logic was. Indeed, *Aristotelian logic*, developed of course by Aristotle and his school, deals with objects, their properties, and quantifiers such as “for all”

and “for some”. A typical line of reasoning (or *syllogism*) in Aristotelean logic goes like this: “All men are mortal. Socrates is a man. Hence, Socrates is mortal”. Aristotelean logic is a subset of mathematical logic, but is not as expressive because it lacks the concept of logical connectives such as and, or, or if-then (although “not” is allowed), and also lacks the concept of a binary relation such as = or <.)

Swapping the order of two quantifiers may or may not make a difference to the truth of a statement. Swapping two “for all” quantifiers is harmless: a statement such as

For all real numbers a , and for all real numbers b , we have $(a+b)^2 = a^2+2ab+b^2$

is logically equivalent to the statement

For all real numbers b , and for all real numbers a , we have $(a+b)^2 = a^2+2ab+b^2$

(why? The reason has nothing to do with whether the identity $(a+b)^2 = a^2+2ab+b^2$ is actually true or not). Similarly, swapping two “there exists” quantifiers has no effect:

There exists a real number a , and there exists a real number b , we have $a^2+b^2 = 0$

is logically equivalent to

There exists a real number b , and there exists a real number a , we have $a^2+b^2 = 0$.

However, swapping a “for all” with a “there exists” makes a lot of difference. Consider the following two statements:

For every integer n , there exists an integer m which is larger than n .
(12.1)

There exists an integer m such that for every integer n , m is larger than n .
(12.2)

Statement 12.1 is obviously true: if your opponent hands you an integer n , you can always find an integer m which is larger than

n . But Statement 12.2 is false: if you choose m first, then you cannot ensure that m is larger than every integer n ; your opponent can easily pick a number n bigger than m to defeat that. The crucial difference between the two statements is that in Statement 12.1, the integer n was chosen *first*, and integer m could then be chosen in a manner depending on n ; but in Statement 12.2, one was forced to choose m first, without knowing in advance what n is going to be. In short, the reason why the order of quantifiers is important is that the inner variables may possibly depend on the outer variables, but not vice versa.

Exercise 12.5.1. What do each of the following statements mean, and which of them are true? Can you find gaming metaphors for each of these statements?

- For every positive number x , and every positive number y , we have $y^2 = x$.
- There exists a positive number x such that for every positive number y , we have $y^2 = x$.
- There exists a positive number x , and there exists a positive number y , such that $y^2 = x$.
- For every positive number y , there exists a positive number x such that $y^2 = x$.
- There exists a positive number y such that for every positive number x , we have $y^2 = x$.

12.6 Some examples of proofs and quantifiers

Here we give some simple examples of proofs involving the “for all” and “there exists” quantifiers. The results themselves are simple, but you should pay attention instead to how the quantifiers are arranged and how the proofs are structured.

Proposition 12.6.1. *For every $\varepsilon > 0$ there exists a $\delta > 0$ such that $2\delta < \varepsilon$.*

Proof. Let $\varepsilon > 0$ be arbitrary. We have to show that there exists a $\delta > 0$ such that $2\delta < \varepsilon$. We only need to pick one such δ ; choosing $\delta := \varepsilon/3$ will work, since one then has $2\delta = 2\varepsilon/3 < \varepsilon$. \square

Notice how ε has to be arbitrary, because we are proving something for *every* ε ; on the other hand, δ can be chosen as you wish, because you only need to show that there *exists* a δ which does what you want. Note also that δ can depend on ε , because the δ -quantifier is nested inside the ε -quantifier. If the quantifiers were reversed, i.e., if you were asked to prove “There exists a $\delta > 0$ such that for every $\varepsilon > 0$, $2\delta < \varepsilon$ ”, then you would have to select δ *first* before being given ε . In this case it is impossible to prove the statement, because it is false (why?).

Normally, when one has to prove a “There exists...” statement, e.g., “Prove that there exists an $\varepsilon > 0$ such that X is true”, one proceeds by selecting ε carefully, and then showing that X is true for that ε . However, this sometimes requires a lot of foresight, and it is legitimate to defer the selection of ε until later in the argument, when it becomes clearer what properties ε needs to satisfy. The only thing to watch out for is to make sure that ε does not depend on any of the bound variables nested inside X . For instance:

Proposition 12.6.2. *There exists an $\varepsilon > 0$ such that $\sin(x) > x/2$ for all $0 < x < \varepsilon$.*

Proof. We pick $\varepsilon > 0$ to be chosen later, and let $0 < x < \varepsilon$. Since the derivative of $\sin(x)$ is $\cos(x)$, we see from the mean-value theorem we have

$$\frac{\sin(x)}{x} = \frac{\sin(x) - \sin(0)}{x - 0} = \cos(y)$$

for some $0 \leq y \leq x$. Thus in order to ensure that $\sin(x) > x/2$, it would suffice to ensure that $\cos(y) > 1/2$. To do this, it would suffice to ensure that $0 \leq y < \pi/3$ (since the cosine function takes the value of 1 at 0, takes the value of $1/2$ at $\pi/3$, and is decreasing in between). Since $0 \leq y \leq x$ and $0 < x < \varepsilon$, we

see that $0 \leq y < \varepsilon$. Thus if we pick $\varepsilon := \pi/3$, then we have $0 \leq y < \pi/3$ as desired, and so we can ensure that $\sin(x) > x/2$ for all $0 < x < \varepsilon$. \square

Note that the value of ε that we picked at the end did not depend on the nested variables x and y . This makes the above argument legitimate. Indeed, we can rearrange it so that we don't have to postpone anything:

Proof. We choose $\varepsilon := \pi/3$; clearly $\varepsilon > 0$. Now we have to show that for all $0 < x < \pi/3$, we have $\sin(x) > x/2$. So let $0 < x < \pi/3$ be arbitrary. By the mean-value theorem we have

$$\frac{\sin(x)}{x} = \frac{\sin(x) - \sin(0)}{x - 0} = \cos(y)$$

for some $0 \leq y \leq x$. Since $0 \leq y \leq x$ and $0 < x < \pi/3$, we have $0 \leq y < \pi/3$. Thus $\cos(y) > \cos(\pi/3) = 1/2$, since \cos is decreasing on the interval $[0, \pi/3]$. Thus we have $\sin(x)/x > 1/2$ and hence $\sin(x) > x/2$ as desired. \square

If we had chosen ε to depend on x and y then the argument would not be valid, because ε is the outer variable and x, y are nested inside it.

12.7 Equality

As mentioned before, one can create statements by starting with expressions (such as $2 \times 3 + 5$) and then asking whether an expression obeys a certain property, or whether two expressions are related by some sort of relation ($=, \leq, \in$, etc.). There are many relations, but the most important one is *equality*, and it is worth spending a little time reviewing this concept.

Equality is a relation linking two objects x, y of the same type T (e.g., two integers, or two matrices, or two vectors, etc.). Given two such objects x and y , the statement $x = y$ may or may not be true; it depends on the value of x and y and also on how equality is defined for the class of objects under consideration. For instance,

for the real numbers the two numbers $0.9999\dots$ and 1 are equal. In modular arithmetic with modulus 10 (in which numbers are considered equal to their remainders modulo 10), the numbers 12 and 2 are considered equal, $12 = 2$, even though this is not the case in ordinary arithmetic.

How equality is defined depends on the class T of objects under consideration, and to some extent is just a matter of definition. However, for the purposes of logic we require that equality obeys the following four *axioms of equality*:

- (Reflexive axiom). Given any object x , we have $x = x$.
- (Symmetry axiom). Given any two objects x and y of the same type, if $x = y$, then $y = x$.
- (Transitive axiom). Given any three objects x , y , z of the same type, if $x = y$ and $y = z$, then $x = z$.
- (Substitution axiom). Given any two objects x and y of the same type, if $x = y$, then $f(x) = f(y)$ for all functions or operations f . Similarly, for any property $P(x)$ depending on x , if $x = y$, then $P(x)$ and $P(y)$ are equivalent statements.

The first three axioms are clear, together, they assert that equality is an *equivalence relation*. To illustrate the substitution we give some examples.

Example 12.7.1. Let x and y be real numbers. If $x = y$, then $2x = 2y$, and $\sin(x) = \sin(y)$. Furthermore, $x + z = y + z$ for any real number z .

Example 12.7.2. Let n and m be integers. If n is odd and $n = m$, then m must also be odd. If we have a third integer k , and we know that $n > k$ and $n = m$, then we also know that $m > k$.

Example 12.7.3. Let x, y, z be real numbers. If we know that $x = \sin(y)$ and $y = z^2$, then (by the substitution axiom) we have $\sin(y) = \sin(z^2)$, and hence (by the transitive axiom) we have $x = \sin(z^2)$.

Thus, from the point of view of logic, we can define equality on a however we please, so long as it obeys the reflexive, symmetry, and transitive axioms, and it is consistent with all other operations on the class of objects under discussion in the sense that the substitution axiom was true for all of those operations. For instance, if we decided one day to modify the integers so that 12 was now equal to 2, one could only do so if one also made sure that 2 was now equal to 12, and that $f(2) = f(12)$ for any operation f on these modified integers. For instance, we now need $2 + 5$ to be equal to $12 + 5$. (In this case, pursuing this line of reasoning will eventually lead to modular arithmetic with modulus 10.)

Exercise 12.7.1. Suppose you have four real numbers a, b, c, d and you know that $a = b$ and $c = d$. Use the above four axioms to deduce that $a + d = b + c$.

Chapter 13

Appendix: the decimal system

In Chapters 2, 4, 5 we painstakingly constructed the basic number systems of mathematics: the natural numbers, integers, rationals, and reals. Natural numbers were simply postulated to exist, and to obey five axioms; the integers then came via (formal) differences of the natural numbers; the rationals then came from (formal) quotients of the integers, and the reals then came from (formal) limits of the rationals.

This is all very well and good, but it does seem somewhat alien to one's prior experience with these numbers. In particular, very little use was made of the *decimal system*, in which the digits 0, 1, 2, 3, 4, 5, 6, 7, 8, 9 are combined to represent these numbers. Indeed, except for a number of examples which were not essential to the main construction, the only decimals we really used were the numbers 0, 1, and 2, and the latter two can be rewritten as 0_{++} and $(0_{++})_{++}$.

The basic reason for this is that *the decimal system itself is not essential to mathematics*. It is very convenient for computations, and we have grown accustomed to it thanks to a thousand years of use, but in the history of mathematics it is actually a comparatively recent invention. Numbers have been around for about ten thousand years (starting from scratch marks on cave walls), but the modern Hindi-Arabic base 10 system for representing numbers only dates from the 11th century or so. Some early civilizations relied on other bases; for instance the Babyloni-

ans used a base 60 system (which still survives in our time system of hours, minutes, and seconds, and in our angular system of degrees, minutes, and seconds). And the ancient Greeks were able to do quite advanced mathematics, despite the fact that the most advanced number representation system available to them was the Roman numeral system I, II, III, IV, \dots , which was horrendous for computations of even two-digit numbers. And of course modern computing relies on binary, hexadecimal, or byte-based (base 256) arithmetic instead of decimal, while analog computers such as the slide rule do not really rely on any number representation system at all. In fact, now that computers can do the menial work of number-crunching, there is very little use for decimals in modern mathematics. Indeed, we rarely use any numbers other than one-digit numbers or one-digit fractions (as well as e, π, i) explicitly in modern mathematical work; any more complicated numbers usually get called more generic names such as n .

Nevertheless, the subject of decimals does deserve an appendix, because it is so integral to the way we use mathematics in our everyday life, and also because we do want to use such notation as $3.14159\dots$ to refer to real numbers, as opposed to the far clunkier “ $LIM_{n \rightarrow \infty} a_n$, where $a_1 = 3.1, a_2 := 3.14, a_3 := 3.141, \dots$ ”.

We begin by reviewing how the decimal system works for the positive integers, and then continue on to the reals. Note that in this discussion we shall freely use all the results from earlier chapters.

13.1 The decimal representation of natural numbers

In this section we will avoid the usual convention of abbreviating $a \times b$ as ab , since this would mean that decimals such as 34 might be misconstrued as 3×4 .

Definition 13.1.1 (Digits). A *digit* is any one of the ten symbols $0, 1, 2, 3, \dots, 9$. We equate these digits with natural numbers by the formulae $0 \equiv 0, 1 \equiv 0++, 2 \equiv 1++,$ etc. all the way up to $9 \equiv 8++$. We also define the number ten by the formula $\text{ten} := 9++$. (We cannot use the decimal notation 10 to denote ten yet,

because that presumes knowledge of the decimal system and would be circular.)

Definition 13.1.2 (Positive integer decimals). A *positive integer decimal* is any string $a_n a_{n-1} \dots a_0$ of digits, where $n \geq 0$ is a natural number, and the first digit a_n is non-zero. Thus, for instance, 3049 is a positive integer decimal, but 0493 or 0 is not. We equate each positive integer decimal with a positive integer by the formula

$$a_n a_{n-1} \dots a_0 \equiv \sum_{i=0}^n a_i \times \text{ten}^i.$$

Remark 13.1.3. Note in particular that this definition implies that

$$10 = 0 \times \text{ten}^0 + 1 \times \text{ten}^1 = \text{ten}$$

and thus we can write ten as the more familiar 10. Also, a single digit integer decimal is exactly equal to that digit itself, e.g., the decimal 3 by the above definition is equal to

$$3 = 3 \times \text{ten}^0 = 3$$

so there is no possibility of confusion between a single digit, and a single digit decimal. (This is a subtle distinction, and not one which is worth losing much sleep over.)

Now we show that this decimal system indeed represents the positive integers. It is clear from the definition that every positive decimal representation gives a positive integer, since the sum consists entirely of natural numbers, and the last term $a_n \text{ten}^n$ is non-zero by definition.

Theorem 13.1.4 (Uniqueness and existence of decimal representations). *Every positive integer m is equal to exactly one positive integer decimal (which is known as the decimal representation of m).*

Proof. We shall use the principle of strong induction (Proposition 2.2.14, with $m_0 := 1$). For any positive integer m , let $P(m)$ denote

the statement “ m is equal to exactly one positive integer decimal”. Suppose we already know $P(m')$ is true for all positive integers $m' < m$; we now wish to prove $P(m)$.

First observe that either $m \geq \text{ten}$ or $m \in \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$. (This is easily proved by ordinary induction.) Suppose first that $m \in \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$. Then m clearly is equal to a positive integer decimal consisting of a single digit, and there is only one single-digit decimal which is equal to m . Furthermore, no decimal consisting of two or more digits can equal m , since if $a_n \dots a_0$ is such a decimal (with $n > 0$) we have

$$a_n \dots a_0 = \sum_{i=0}^n a_i \times \text{ten}^i \geq a_n \times \text{ten}^i \geq \text{ten} > m.$$

Now suppose that $m \geq \text{ten}$. Then by the Euclidean algorithm (Proposition 2.3.9) we can write

$$m = s \times \text{ten} + r$$

where s is a positive integer, and $r \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$. Since

$$s < s \times \text{ten} \leq a \times \text{ten} + r = m$$

we can use the strong induction hypothesis and conclude that $P(s)$ is true. In particular, s has a decimal representation

$$s = b_p \dots b_0 = \sum_{i=0}^p b_i \times \text{ten}^i.$$

Multiplying by ten, we see that

$$s \times \text{ten} = \sum_{i=0}^p b_i \times \text{ten}^{i+1} = b_p \dots b_0 0,$$

and then adding r we see that

$$m = s \times \text{ten} + r = \sum_{i=0}^p b_i \times \text{ten}^{i+1} + r = b_p \dots b_0 r.$$

Thus m has at least one decimal representation. Now we need to show that m has at most one decimal representation. Suppose for contradiction that we have at least two different representations

$$m = a_n \dots a_0 = a'_{n'} \dots a'_0.$$

First observe by the previous computation that

$$a_n \dots a_0 = (a_n \dots a_1) \times \text{ten} + a_0$$

and

$$a'_{n'} \dots a'_0 = (a'_{n'} \dots a'_1) \times \text{ten} + a'_0$$

and so after some algebra we obtain

$$a_0 - a'_0 = (a_n \dots a_1 - a'_{n'} \dots a'_1) \times \text{ten}.$$

The right-hand side is a multiple of ten, while the left-hand side lies strictly between $-\text{ten}$ and $+\text{ten}$. Thus both sides must be equal to 0. This means that $a_0 = a'_0$ and $a_n \dots a_1 = a'_{n'} \dots a'_1$. But by the previous arguments, we know that $a_n \dots a_1$ is a smaller integer than $a_n \dots a_0$. Thus by the strong induction hypothesis, the number $a_n \dots a_0$ has only one decimal representation, which means that n' must equal n and a'_i must equal a_i for all $i = 1, \dots, n$. Thus the decimals $a_n \dots a_0$ and $a'_{n'} \dots a'_0$ are in fact identical, contradicting the assumption that they were different. \square

Once one has decimal representation, one can then derive the usual laws of long addition and long multiplication to connect the decimal representation of $x + y$ or $x \times y$ to that of x or y (Exercise 13.1.1).

Once one has decimal representation of positive integers, one can of course represent negative integers decimally as well by using the $-$ sign. Finally, we let 0 be a decimal as well. This gives decimal representations of all integers. Every rational is then the ratio of two decimals, e.g., $335/113$ or $-1/2$ (with the denominator required to be non-zero, of course), though of course there may

be more than one way to represent a rational as such a ratio, e.g., $6/4 = 3/2$.

Since ten = 10, we will now use 10 instead of ten throughout, as is customary.

Exercise 13.1.1. The purpose of this exercise is to demonstrate that the procedure of long addition taught to you in elementary school is actually valid. Let $A = a_n \dots a_0$ and $B = b_m \dots b_0$ be positive integer decimals. Let us adopt the convention that $a_i = 0$ when $i > n$, and $b_i = 0$ when $i > m$; for instance, if $A = 372$, then $a_0 = 2$, $a_1 = 7$, $a_2 = 3$, $a_3 = 0$, $a_4 = 0$, and so forth. Define the numbers c_0, c_1, \dots and $\varepsilon_0, \varepsilon_1, \dots$ recursively by the following *long addition algorithm*.

- We set $\varepsilon_0 := 0$.
- Now suppose that ε_i has already been defined for some $i \geq 0$. If $a_i + b_i + \varepsilon_i < 10$, we set $c_i := a_i + b_i + \varepsilon_i$ and $\varepsilon_{i+1} := 0$; otherwise, if $a_i + b_i + \varepsilon_i \geq 10$, we set $c_i := a_i + b_i + \varepsilon_i - 10$ and $\varepsilon_{i+1} = 1$. (The number ε_{i+1} is the “carry digit” from the i^{th} decimal place to the $(i + 1)^{\text{th}}$ decimal place.)

Prove that the numbers c_0, c_1, \dots are all digits, and that there exists an l such that $c_l \neq 0$ and $c_i = 0$ for all $i > l$. Then show that $c_l c_{l-1} \dots c_1 c_0$ is the decimal representation of $A + B$.

Note that one could in fact use this algorithm to *define* addition, but it would look extremely complicated, and to prove even such simple facts as $(a + b) + c = a + (b + c)$ would be rather difficult. This is one of the reasons why we have avoided the decimal system in our construction of the natural numbers. The procedure for long multiplication (or long subtraction, or long division) is even worse to lay out rigorously; we will not do so here.

13.2 The decimal representation of real numbers

We need a new symbol: the *decimal point* “.”.

Definition 13.2.1 (Real decimals). A *real decimal* is any sequence of digits, and a decimal point, arranged as

$$\pm a_n \dots a_0 . a_{-1} a_{-2} \dots$$

which is finite to the left of the decimal point (so n is a natural number), but infinite to the right of the decimal point, where \pm is either $+$ or $-$, and $a_n \dots a_0$ is a natural number decimal (i.e., either a positive integer decimal, or 0). This decimal is equated to the real number

$$\pm a_n \dots a_0 . a_{-1} a_{-2} \dots \equiv \pm 1 \times \sum_{i=-\infty}^N a_i \times 10^i.$$

The series is always convergent (Exercise 13.2.1). Next, we show that every real number has at least one decimal representation:

Theorem 13.2.2 (Existence of decimal representations). *Every real number x has at least one decimal representation $\pm a_n \dots a_0 . a_{-1} a_{-2} \dots$.*

Proof. We first note that $x = 0$ has the decimal representation $0.000 \dots$. Also, once we find a decimal representation for x , we automatically get a decimal representation for $-x$ by changing the sign \pm . Thus it suffices to prove the theorem for real numbers x (by Proposition 5.4.4).

Let $n \geq 0$ be any natural number. From the Archimedean property (Corollary 5.4.13) we know that there is a natural number M such that $M \times 10^{-n} > x$. Since $0 \times 10^{-n} \leq x$, we thus see that there must exist a natural number s_n such that $s_n \times 10^{-n} \leq x$ and $s_n + 1 \times 10^{-n} > x$. (If no such natural number existed, one could use induction to conclude that $s \times 10^{-n} \leq x$ for all natural numbers s , contradicting the Archimedean property, Corollary 5.4.13.)

Now consider the sequence s_0, s_1, s_2, \dots . Since we have

$$s_n \times 10^{-n} \leq x < (s_n + 1) \times 10^{-n}$$

we thus have

$$(10 \times s_n) \times 10^{-(n+1)} \leq x < (10 \times s_n + 10) \times 10^{-(n+1)}.$$

On the other hand, we have

$$s_{n+1} \times 10^{-(n+1)} \leq x < (s_{n+1} + 1) \times 10^{-(n+1)}$$

and hence we have

$$10 \times s_n < s_{n+1} + 1 \text{ and } s_{n+1} \leq 10 \times s_n + 10.$$

From these two inequalities we see that we have

$$10 \times s_n \leq s_{n+1} \leq 10 \times s_n + 9$$

and hence we can find a digit a_{n+1} such that

$$s_{n+1} = 10 \times s_n + a_{n+1}$$

and hence

$$s_{n+1} \times 10^{-(n+1)} = s_n \times 10^{-n} + a_{n+1} \times 10^{-(n+1)}.$$

From this identity and induction, we can obtain the formula

$$s_n \times 10^{-n} = s_0 + \sum_{i=0}^n a_i \times 10^{-i}.$$

Now we take limits of both sides (using Exercise 13.2.1) to obtain

$$\lim_{n \rightarrow \infty} s_n \times 10^{-n} = s_0 + \sum_{i=0}^{\infty} a_i \times 10^{-i}.$$

On the other hand, we have

$$x - 10^{-n} \leq s_n \times 10^{-n} \leq x$$

for all n , so by the squeeze test (Corollary 6.3.14) we have

$$\lim_{n \rightarrow \infty} s_n \times 10^{-n} = x.$$

Thus we have

$$x = s_0 + \sum_{i=0}^{\infty} a_i \times 10^{-i}.$$

Since s_0 already has a positive integer decimal representation by Theorem 1, we thus see that x has a decimal representation as desired. \square

There is however one slight flaw with the decimal system: it is possible for one real number to have two decimal representations.

Proposition 13.2.3 (Failure of uniqueness of decimal representations). *The number 1 has two different decimal representations: 1.000... and 0.999....*

Proof. The representation $1 = 1.000\dots$ is clear. Now let's compute $0.999\dots$. By definition, this is the limit of the Cauchy sequence

$$0.9, 0.99, 0.999, 0.9999, \dots$$

But this sequence has 1 as a formal limit by Proposition 5.2.8. \square

It turns out that these are the only two decimal representations of 1 (Exercise 13.2.2). In fact, as it turns out, all real numbers have either one or two decimal representations - two if the real is a terminating decimal, and one otherwise (Exercise 13.2.3).

Exercise 13.2.1. If $a_n \dots a_0.a_{-1}a_{-2} \dots$ is a real decimal, show that the series $\sum_{i=-\infty}^N a_i \times 10^i$ is absolutely convergent.

Exercise 13.2.2. Show that the only decimal representations $1 = \pm a_n \dots a_0.a_{-1}a_{-2} \dots$ of 1 are $1 = 1.000\dots$ and $1 = 0.999\dots$

Exercise 13.2.3. A real number x is said to be a *terminating decimal* if we have $x = n/10^{-m}$ for some integers n, m . Show that if x is a terminating decimal, then x has exactly two decimal representations, while if x is not a terminating decimal, then x has exactly one decimal representation.

Exercise 13.2.4. Rewrite the proof of Corollary 8.3.4 using the decimal system.

Index

- ++ (increment), 8, 56
 - on integers, 88
- + C , 345
- α -length, 336
- ε -adherent, 161, 245
 - continually ε -adherent, 161
- ε -close
 - eventually ε -close sequences, 116
 - functions and numbers, 254
 - locally ε -close functions and numbers, 255
 - rational numbers, 100
 - reals, 146
 - sequences, 116, 148
- ε -steady, 111, 147
 - eventually ε -steady, 112, 147
- π , 510, 512
- σ -algebra, 580, 598

- a posteriori*, 20, 303
- a priori*, 20, 303
- Abel's theorem, 487
- absolute convergence
 - for series, 192, 217
 - test, 192
- absolute value
 - for complex numbers, 502
 - for rationals, 100
 - for reals, 129
- absolutely integrable, 620
- absorption laws, 52
- abstraction, 24-25, 391
- addition
 - long, 387
 - of cardinals, 82
 - of functions, 253
 - of complex numbers, 119
 - of integers, 87, 88
 - of natural numbers, 27
 - of rational numbers, 94, 95
 - of reals, 119-121
- (countably) additive measure, 580, 596
- adherent point
 - infinite, 288
 - of sequences: *see* limit point
 - of sequences
 - of sets: 246, 404, 438
- alternating series test, 193
- ambient space, 408
- analysis, 1
- and: *see* conjunction
- antiderivative, 343
- approximation to the identity, 468, 472, 527
- Archimedean property, 132

- arctangent: *see* trigonometric functions
- Aristotlean logic, 375
- associativity
- of composition, 60
 - of addition in \mathbf{C} , 499
 - of addition in \mathbf{N} , 29
 - of addition in vector spaces, 538
 - of multiplication in \mathbf{N} , 34
 - of scalar multiplication, 538
 - see also*: ring, field, laws of algebra
- asymptotic discontinuity, 269
- Axiom(s)
- in mathematics, 24-25
 - of choice, 40, 75, 229
 - of comprehension: *see* Axiom of universal specification
 - of countable choice, 230
 - of equality, 380
 - of foundation: *see* Axiom of regularity
 - of induction: *see* principle of mathematical induction
 - of infinity, 50
 - of natural numbers: *see* Peano axioms
 - of pairwise union, 42
 - of power set, 66
 - of reflexivity, 380
 - of regularity, 54
 - of replacement, 49
 - of separation, 45
 - of set theory, 38, 40-42, 45, 49-50, 54, 66
 - of singleton sets and pair sets, 41
 - of specification, 45
 - of symmetry, 380
 - of substitution, 57, 380
 - of the empty set, 40
 - of transitivity, 380
 - of universal specification, 52
 - of union, 67
- ball, 402
- Banach-Tarski paradox, 576, 594
- base of the natural logarithm: *see e*
- basis
- standard basis of row vectors, 539
- bijection, 62
- binomial formula, 189
- Bolzano-Weierstrass theorem, 174
- Boolean algebra, 47, 580, 595
- Boolean logic, 369
- Borel property, 579, 599
- Borel-Cantelli lemma, 618
- bound variable, 180, 371, 396
- boundary (point), 403, 437
- bounded
- from above and below, 270
 - function, 270, 454
 - interval, 245
 - sequence, 114, 151

- sequence away from zero,
 - 124, 128
- set, 249, 415
- C, C^0, C^1, C^2, C^k , 560
- cancellation law
 - of addition in \mathbf{N} , 29
 - of multiplication in \mathbf{N} , 34
 - of multiplication in \mathbf{Z} , 92
 - of multiplication in \mathbf{R} , 127
- Cantor's theorem, 224
- cardinality
 - arithmetic, 82
 - of finite sets, 80
 - uniqueness of, 80
- Cartesian product, 71
 - infinite, 229
- Cauchy criterion, 197
- Cauchy sequence, 112, 146, 411
- Cauchy-Schwarz inequality, 401, 520
- chain: *see* totally ordered set
- chain rule, 295
 - in higher dimensions, 556, 559
- change of variables formula, 348-350
- character, 522
- characteristic function, 606
- choice
 - single, 40
 - finite, 74
 - countable, 230
 - arbitrary, 229
- closed
 - box, 584
 - interval, 244
 - set, 405, 438
- Clairaut's theorem: *see* interchanging derivatives with derivatives
- closure, 246, 404, 438
- cluster point: *see* limit point
- cocountable topology, 440
- coefficient, 477
- cofinite topology, 440
- column vector, 538
- common refinement, 313
- commutativity
 - of addition in \mathbf{C} , 499
 - of addition in \mathbf{N} , 29
 - of addition in vector spaces, 538
 - of convolution, 470, 526
 - of multiplication in \mathbf{N} , 34
 - see also*: ring, field, laws of algebra
- compactness, 415, 439
 - compact support, 469
- comparison principle (or test)
 - for finite series, 181
 - for infinite series, 196
 - for sequences, 167
- completeness
 - of the space of continuous functions, 457
 - of metric spaces, 412
 - of the reals, 168
- completion of a metric space, 414
- complex numbers \mathbf{C} , 499
 - complex conjugation, 501
- composition of functions, 59

- conjunction (and), 256
- connectedness, 309, 433
 - connected component, 435
- constant
 - function, 58, 314
 - sequence, 171
- continuity, 262, 422, 438
 - and compactness, 429
 - and connectedness, 434
 - and convergence, 256, 423
- continuum, 243
 - hypothesis, 227
- contraction, 562
 - mapping theorem, 563
- contrapositive, 364
- convergence
 - in L^2 , 521
 - of a function at a point, 255, 444
 - of sequences, 148, 396, 437
 - of series, 190
 - pointwise: *see* pointwise convergence
 - uniform: *see* uniform convergence
- converse, 364
- convolution, 470, 491, 526
- Corollary, 28
- coset, 591
- cosine: *see* trigonometric functions
- cotangent: *see* trigonometric functions
- countability, 208
 - of the integers, 212
 - of the rationals, 214
- cover, 582
 - see also*: open cover
- critical point, 576
- de Moivre identities, 511
- de Morgan laws, 47
- decimal
 - negative integer, 386
 - non-uniqueness of representation, 390
 - point, 387
 - positive integer, 384
 - real, 388
- degree, 468
- denumerable: *see* countable
- dense, 468
- derivative, 290
 - directional, 548
 - in higher dimensions, 546, 548, 550, 555
 - partial, 550
 - matrix, 555
 - total, 546, 548
 - uniqueness, 547
- difference rule, 294
- difference set, 47
- differential matrix: *see* derivative matrix
- differentiability
 - at a point, 290
 - continuous, 560
 - directional, 548
 - in higher dimensions, 546
 - infinite, 481
 - k -fold, 481, 560
- digit, 383
- dilation, 540

- diophantine, 619
- Dirac delta function, 469
- direct sum
 - of functions, 76, 425
- discontinuity: *see* singularity
- discrete metric, 395
- disjoint sets, 47
- disjunction (or), 356
 - inclusive vs. exclusive, 356
- distance
 - in \mathbf{Q} , 100
 - in \mathbf{R} , 146, 391
- distributive law
 - for natural numbers, 34
 - for complex numbers, 500
 - see also*: laws of algebra
- divergence
 - of series, 3, 190
 - of sequences, 4
 - see also*: convergence
- divisibility, 238
- division
 - by zero, 3
 - formal ($//$), 94
 - of functions, 253
 - of rationals, 97
- domain, 55
- dominate: *see* majorize
 - dominated convergence: *see*
 - Lebesgue dominated convergence theorem
- doubly infinite, 245
- dummy variable: *see* bound variable
- e , 494
- Egoroff's theorem, 620
- empty
 - Cartesian product, 64
 - function, 59
 - sequence, 64
 - series, 185
 - set, 38, 580, 583
- equality, 379
 - for functions, 58
 - for sets, 39
 - of cardinality, 79
- equivalence
 - of sequences, 177, 283
 - relation, 380
- error-correcting codes, 394
- Euclidean algorithm, 35
- Euclidean metric, 393
- Euclidean space, 393
- Euler's formula, 507, 510
- Euler's number: *see* e
- exponential function, 493, 504
- exponentiation
 - of cardinals, 82
 - with base and exponent in \mathbf{N} , 36
 - with base in \mathbf{Q} and exponent in \mathbf{Z} , 102, 103
 - with base in \mathbf{R} and exponent in \mathbf{Z} , 140, 141
 - with base in \mathbf{R}^+ and exponent in \mathbf{Q} , 143
 - with base in \mathbf{R}^+ and exponent in \mathbf{R} , 177
- expression, 355
- extended real number system
 - \mathbf{R}^* , 138, 155
- extremum: *see* maximum, min-

- imum
- exterior (point), 403, 437
- factorial, 189
- family, 68
- Fatou's lemma, 617
- Fejér kernel, 528
- field, 97
 - ordered, 99
- finite intersection property, 418
- finite set, 81
- fixed point theorem, 277, 563
- forward image: *see* image
- Fourier
 - coefficients, 524
 - inversion formula, 524
 - series, 524
 - series for arbitrary periods, 535
 - theorem, 530
 - transform, 524
- fractional part, 516
- free variable, 370
- frequency, 523
- Fubini's theorem, 628
 - for finite series, 188
 - for infinite series, 217
 - see also*: interchanging integrals/sums with integrals/sums
- function, 55
 - implicit definition, 57
- fundamental theorems of calculus, 340, 343
- geometric series, 196, 197
 - formula, 197, 200, 463
- geodesic, 395
- gradient, 554
- graph, 58, 76, 252, 571
- greatest lower bound: *see* least upper bound
- hairy ball theorem, 563
- half-infinite, 245
- half-open, 244
- half-space, 595
- harmonic series, 199
- Hausdorff space, 437, 439
- Hausdorff maximality principle, 241
- Heine-Borel theorem, 416
 - for the real line, 249
- Hermitian, 519
- homogeneity, 520, 539
- hypersurface, 572
- identity map (or operator), 64, 540
- if: *see* implication
- iff (if and only if), 30
- ill-defined, 353
- image
 - of sets, 64
 - inverse image, 65
- imaginary, 501
- implication (if), 360
- implicit differentiation, 573
- implicit function theorem, 572
- improper integral, 320
- inclusion map, 64
- inconsistent, 506
- index of summation: *see* dummy variable

- index set, 68
- indicator function: *see* characteristic function
- induced
 - metric, 393, 408
 - topology, 408, 438
- induction: *see* Principle of mathematical induction
- infinite
 - interval, 245
 - set, 80
- infimum: *see* supremum
- injection: *see* one-to-one function
- inner product, 518
- integer part, 104, 133, 516
- integers \mathbf{Z}
 - definition 86
 - identification with rationals, 95
 - interspersing with rationals, 104
- integral test, 334
- integration
 - by parts, 345-347, 488
 - laws, 317, 323
 - piecewise constant, 315, 317
 - Riemann: *see* Riemann integral
- interchanging
 - derivatives with derivatives, 10, 560
 - finite sums with finite sums, 187, 188
 - integrals with integrals, 7, 628
- limits with derivatives, 9, 464
- limits with integrals, 9, 462
- limits with length, 12
- limits with limits, 8, 9, 453
- limits with sums, 620
- sums with derivatives, 466, 479
- sums with integrals, 463, 479, 615, 616
- sums with sums, 6, 217
- interior (point), 403, 437
- intermediate value theorem, 275, 434
- intersection
 - pairwise, 46
- interval, 244
- intrinsic, 415
- inverse
 - function theorem, 303, 567
 - image, 65
 - in logic, 364
 - of functions, 63
- invertible function: *see* bijection
 - local, 566
- involution, 501
- irrationality, 138
 - existence of, 105, 138
 - of $\sqrt{2}$, 105
 - need for, 109
- isolated point, 248
- isometry, 408

- jump discontinuity, 269
- l^1 , l^2 , l^∞ , L^1 , L^2 , L^∞ , 393-395, 520, 621
 - equivalence of in finite dimensions, 398
 - see also*: absolutely integrable
 - see also*: supremum as norm
- L'Hôpital's rule, 11, 305, 306
- label, 68
- laws of algebra
 - for complex numbers, 500
 - for integers, 90
 - for rationals, 96
 - for reals, 123
- laws of arithmetic: *see* laws of algebra
- laws of exponentiation, 102, 103, 141, 144, 178, 494
- least upper bound, 135
 - least upper bound property, 136, 159
 - see also* supremum
- Lebesgue dominated convergence theorem, 622
- Lebesgue integral
 - of absolutely integrable functions, 621
 - of nonnegative functions, 612
 - of simple functions, 607
 - upper and lower, 623
 - vs. the Riemann integral, 626
- Lebesgue measurable, 594
- Lebesgue measure, 581
 - motivation of, 579-581
- Lebesgue monotone convergence theorem, 613
- Leibnitz rule, 294, 558
- Lemma, 28
- length of interval, 310
- limit
 - at infinity, 288
 - formal (LIM), 119, 414
 - laws, 151, 257, 503
 - left and right, 267
 - limiting values of functions, 5, 255, 444
 - of sequences, 149
 - pointwise, 447
 - uniform, *see* uniform limit
 - uniqueness of, 148, 257, 399, 445
- limit inferior, *see* limit superior
- limit point
 - of sequences, 161, 411
 - of sets, 248
- limit superior, 163
- linear combination, 539
- linearity
 - approximate, 545
 - of convolution, 471, 527
 - of finite series, 186
 - of limits, 151
 - of infinite series, 194
 - of inner product, 519
 - of integration, 318, 609, 612
 - of transformations, 539

- Lipschitz constant, 300
- Lipschitz continuous, 300
- logarithm (natural), 495
 - power series of, 464, 495
- logical connective, 356
- lower bound: *see* upper bound

- majorize, 319, 611
- manifold, 576
- map: *see* function
- matrix, 540
 - identification with linear transformations, 541-544
- maximum, 233
 - local, 297
 - global, 298
 - of functions, 253, 272
 - principle, 272, 430
- mean value theorem, 299
- measurability
 - for functions, 601, 603
 - for sets, 590
 - motivation of, 579
 - see also*: Lebesgue measure, outer measure
- meta-proof, 141
- metric, 392
 - ball: *see* ball
 - on \mathbf{C} , 503
 - on \mathbf{R} , 393
 - space, 392
 - see also*: distance
- minimum, 233
 - local, 297
 - local, 298
 - of a set of natural numbers, 210
 - of functions, 253, 272
- minorize: *see* majorize
- monomial, 522
- monotone (increasing or decreasing)
 - convergence: *see* Lebesgue monotone convergence theorem
 - function, 277, 338
 - measure, 580, 584
 - sequence, 160
- morphism: *see* function
- moving bump example, 449, 617
- multiplication
 - of cardinals, 82
 - of complex numbers, 500
 - of functions, 253
 - of integers, 87
 - of matrices, 540
 - of natural numbers, 33
 - of rationals, 94,95
 - of reals, 121

- Natural numbers \mathbf{N}
 - are infinite, 80
 - axioms: *see* Peano axioms
 - identification with integers, 88
 - informal definition, 17
 - in set theory: *see* Axiom of infinity
 - uniqueness of, 77
- negation
 - in logic, 357

- of extended reals, 155
 - of complex numbers, 499
 - of integers, 89
 - of rationals, 95
 - of reals, 122
- negative: *see* negation, positive
- neighbourhood, 437
- Newton's approximation, 293, 548
- non-constructive, 230
- non-degenerate, 520
- nowhere differentiable function, 467, 512
- objects, 38
 - primitive, 53
- one-to-one function, 61
- one-to-one correspondence: *see* bijection
- onto, 61
- open
 - box, 582
 - cover, 416
 - interval, 244
 - set, 405
- or: *see* disjunction
- order ideal, 238
- order topology, 440
- ordered pair, 71
 - construction of, 75
- ordered n -tuple, 72
- ordering
 - lexicographical, 239
 - of cardinals, 227
 - of orderings, 240
 - of partitions, 312
 - of sets, 233
 - of the extended reals, 155
 - of the integers, 92
 - of the natural numbers, 31
 - of the rationals, 98
 - of the reals, 130
- orthogonality, 520
- orthonormal, 523
- oscillatory discontinuity, 269
- outer measure, 583
 - non-additivity of, 591, 593
- pair set, 41
- partial function, 70
- partially ordered set, 45, 232
- partial sum, 190
- Parseval identity, 535
 - see also*: Plancherel formula
- partition, 310
- path-connected, 435
- Peano axioms, 18-21, 23
- perfect matching: *see* bijection
- periodic, 515
 - extension, 516
- piecewise
 - constant, 314
 - constant Riemann-Stieltjes integral, 337
 - continuous, 332
- Plancherel formula (or theorem), 524, 532
- pointwise convergence, 447
 - of series, 459
 - topology of, 458

- polar representation, 511
- polynomial, 267, 468
 - and convolution, 471
 - approximation by, 468, 470
- positive
 - complex number, 502, 506
 - integer, 90
 - inner product, 519
 - measure, 580, 584
 - natural number, 30
 - rational, 98
 - real, 129
- power series, 479
 - formal, 477
 - multiplication of, 490
 - uniqueness of, 484
- power set, 67
- pre-image: *see* inverse image
- principle of infinite descent, 107
- principle of mathematical induction, 21
 - backwards induction, 33
 - strong induction, 32, 234
 - transfinite, 237
- product rule, *see* Leibnitz rule
- product topology, 458
- projection, 540
- proof
 - by contradiction, 354, 365
 - abstract examples, 366-369, 377-379
- proper subset, 44
- property, 356
- Proposition, 28
- propositional logic, 369
- Pythagoras' theorem, 520
- quantifier, 371
 - existential (for some), 373
 - negation of, 374
 - nested, 374
 - universal (for all), 372
- Quotient: *see* division
- Quotient rule, 295, 559
- radius of convergence, 478
- range, 55
- ratio test, 206
- rational numbers **Q**
 - definition, 94
 - identification with reals, 122
 - interspersing with rationals, 104
 - interspersing with reals, 133
- real analytic, 481
- real numbers **R**
 - are uncountable: *see* uncountability of the reals
 - definition, 118
- real part, 501
- real-valued, 459
- rearrangement
 - of absolutely convergent series, 202
 - of divergent series, 203, 222
 - of finite series, 185
 - of non-negative series, 200
- reciprocal

- of complex numbers, 502
 - of rationals, 96
 - of reals, 126
- recursive definitions, 26,77
- reductio ad absurdum*: *see* proof by contradiction
- relation, 256
- relative topology: *see* induced topology
- removable discontinuity: *see* removable singularity
- removable singularity, 260, 269
- restriction of functions, 251
- Riemann hypothesis, 200
- Riemann integrability, 320
 - closure properties, 323-326
 - failure of, 335
 - of bounded continuous functions, 331
 - of continuous functions on compacta, 330
 - of monotone functions, 333
 - of piecewise continuous bounded functions, 332
 - of uniformly continuous functions, 329
- Riemann integral 320
 - upper and lower, 319
- Riemann sums (upper and lower), 321
- Riemann zeta function, 199
- Riemann-Stieltjes integral, 339
- ring, 90
 - commutative, 90, 500
- Rolle's theorem, 299
- root, 141
 - mean square: *see* L^2 test, 204
- row vector, 537
- Russell's paradox, 52
- scalar multiplication, 537
 - of functions, 253
- Schröder-Bernstein theorem, 227
- sequence, 110
 - finite, 74
- series
 - finite, 179, 182
 - formal infinite, 189
 - laws, 194, 220
 - of functions, 459
 - on arbitrary sets, 220
 - on countable sets, 216
 - vs. sum, 180
- set
 - axioms: *see* axioms of set theory
 - informal definition, 38
 - signum function, 259
 - simple function, 605
 - sine: *see* trigonometric functions
 - singleton set, 41
 - singularity, 270
 - space, 392
 - statement, 352
 - sub-additive measure, 580, 584
 - subset, 44
 - subsequence, 173, 410
 - substitution
 - see also*: rearrangement, 185

- subtraction
 - formal ($--$), 86
 - of functions, 253
 - of integers, 91
- sum rule, 294
- summation by parts, 487
- sup norm: *see* supremum as norm
- support, 469
- supremum (and infimum)
 - as metric, 395
 - as norm, 395, 460
 - of a set of extended reals, 156, 157
 - of a set of reals, 137, 139
 - of sequences of reals, 158
- square root, 56
- square wave, 516, 522
- Squeeze test
 - for sequences, 167
- Stone-Weierstrass theorem, 475, 526
- strict upper bound, 235
- surjection: *see* onto
- taxicab metric, 394
- tangent: *see* trigonometric function
- Taylor series, 483
- Taylor's formula: *see* Taylor series
- telescoping series, 195
- ten, 383
- Theorem, 28
- topological space, 436
- totally bounded, 420
- totally ordered set, 45, 233
- transformation: *see* function
- translation, 540
 - invariance, 581, 584, 595
- transpose, 538
- triangle inequality
 - in Euclidean space, 401
 - in inner product spaces, 520
 - in metric spaces, 392
 - in \mathbf{C} , 502
 - in \mathbf{R} , 100
 - for finite series, 181, 186
 - for integrals, 621
- trichotomy of order
 - of extended reals, 156
 - for natural numbers, 31
 - for integers, 92
 - for rationals, 98
 - for reals, 130
- trigonometric functions, 507, 511
 - and Fourier series, 533
- trigonometric polynomials, 522
 - power series, 508, 512
- trivial topology, 439
- two-to-one function, 61
- uncountability, 208
 - of the reals, 225
- undecidable, 228
- uniform continuity, 282, 430
- uniform convergence, 450
 - and anti-derivatives, 465
 - and derivatives, 454
 - and integrals, 462
 - and limits, 453

- and radius of convergence, 479
 - as a metric, 456, 517
 - of series, 459
- uniform limit, 450
 - of bounded functions, 454
 - of continuous functions, 453
 - and Riemann integration, 462
- union, 67
 - pairwise, 42
- universal set, 53
- upper bound,
 - of a set of reals, 134
 - of a partially ordered set, 235
 - see also*: least upper bound
- variable, 370
- vector space, 538
- vertical line test, 55, 76, 572
- volume, 582
- Weierstrass approximation theorem, 468, 473-474, 525
- Weierstrass example: *see* nowhere differentiable function
- Weierstrass M -test, 460
- well-defined, 353
- well-ordered sets, 234
- well ordering principle
 - for natural numbers, 210
 - for arbitrary sets, 241
- Zermelo-Fraenkel(-Choice) ax-
 - ioms, 69
 - see also* axioms of set theory
- zero test
 - for sequences, 168
 - for series, 191
- Zorn's lemma, 237