

Patterns of Structure and Argument:
a guide to mathematical thinking

February 7, 2012

Contents

1	Theory & Proof:	
	Why do I need to know this?	1
1.1	What is mathematics, and what is it good for?	1
1.2	Exercises	4
1.3	Back to the bees! Analyzing the honeycomb.	4
1.4	Exercises	5
1.5	Examining our assumptions.	5
1.6	Extra credit exercises!	6
1.7	Mathematics is abstract by nature.	6
1.8	The dance of abstraction and concreteness.	9
1.9	A note on how this course will proceed.	11
2	Motivation: Area Revisited	13
2.1	What is Area?	13
2.2	Exercises	16
2.3	What have we left out?	18
3	Motivation: Solution Sets	21
3.1	What does solving mean exactly?	21
3.2	The process of solving an equation	25
3.3	Exercises	30
3.4	A More Complicated Equation	31
3.5	Exercises	32
4	Sets, Logic, & Proof	35
4.1	Open Propositions Define Sets	35
4.2	Exercises	41
4.3	Logical Equivalence	41

4.3.1	Contrapositives	42
4.3.2	Rephrasing a conditional as a disjunction	43
4.3.3	Negation: Don't just say no!	44
4.4	Quantified Propositions Define Relationships	47
4.4.1	Quantifiers	47
4.4.2	Relationships Defined by Quantified Propositions	50
4.5	Exercises	51
4.6	The Structure of a Proof	54
4.6.1	From Universal to Specific	55
4.6.2	Generality through Arbitrariness	55
4.6.3	Vacuity	56
4.6.4	Division into cases	56
4.6.5	Proof by Contradiction	59
4.7	Exercises	61
5	The Real & Natural Number Systems	63
5.1	Addition and Multiplication of Real Numbers	65
5.2	Exercises	67
5.3	The Ordering of the Real Numbers	70
5.4	Exercises	72
5.5	The Definition of the Natural Number System	74
5.6	Proof by Induction	75
5.7	Recursive Definition	77
5.8	Proving Inequalities by Induction	80
5.9	Exercises.	81
5.10	Getting the Most out of Induction	83
5.10.1	Extending to Initial Numbers Other Than Zero	83
5.10.2	Complete Induction: Two Approaches	87
5.11	Exercises	91
6	Relations	95
6.1	Ordered pairs	96
6.2	Exercises	96
6.3	The Cartesian Product Operation	97
6.4	Exercises	97
6.5	The Definition of a Relation	98
6.6	Operations on Relations: Inverse and Composition	99
6.7	Exercises	100

6.8	Functions	101
6.8.1	The Definition of a Function	101
6.8.2	Some special types of functions	103
6.8.3	Surjectivity, Injectivity, and Inverse Functions	104
6.8.4	Exercises	106
6.8.5	The Induced Function on Power Sets	108
6.8.6	Exercises	109
6.9	Properties of Relations: a Compendium	109
6.9.1	Exercises	110
6.10	Order Relations	111
6.10.1	Weak Orders Versus Strong Orders	111
6.10.2	Partial Orders Versus Total Orders	112
6.10.3	The Dictionary Order on a Cartesian Product	112
6.10.4	Exercises	112
6.10.5	Maxima and Minima; Bounds; Suprema and Infima	114
6.10.6	Exercises	117
6.11	Equivalence Relations	120
6.11.1	Equivalence Classes and Partitions	120
6.11.2	Exercises	122

Chapter 1

Theory & Proof: Why do I need to know this?

In which we attempt to answer the question posed by the chapter title and, in doing so, provide some useful background on the nature of mathematics.

You may wonder why it is necessary to carefully study mathematical theory. Why study sets, for example? Aren't they just collections of things? Why learn to prove theorems? Other people have proven the theorems we'll use, and why bother to prove them anyway? Can't we tell by experience what works and what doesn't? We've solved equations, computed derivatives, and worked out applied problems without learning the proofs of our methods, and things have gone just fine!

If you have questions and doubts similar to these, they deserve an answer. This introduction is an attempt to give one. I hope it will also help you see that you have used logical proof and theory more than you realize; you just weren't aware of it.

1.1 What is mathematics, and what is it good for?

Presumably, since you are likely a mathematics major or minor, you think mathematics is worthwhile. But let's examine more carefully what mathematics is, and why it is valuable. These are big questions, of course, and this short discussion will barely begin to answer them. It is doubtful they can be fully answered at all. However, I think we can gain some useful and motivational insights from a short examination.

Mathematics is often described as the study of patterns. What is a pattern? A pattern is something that is regular and predictable, and therefore recognizable. The

relationships among its parts can be described with precision; it is these relationships that give the pattern its *structure*.

Consider, for example, a honeycomb made by bees. (If you haven't seen one, just google "honeycomb" images.) It has a regular pattern in that it fills a flat panel with a grid of hexagonal cylinders, each of which is very close to regular. By *hexagonal* we mean six-sided; by *regular* we mean that the sides are the same length (in cross-section) and the corner angles have equal measure. (In other words, the sides and angles are all congruent.)

In order to fully explain our description of the honeycomb, we'd have to precisely specify or define many concepts of plane geometry, such as lines, segments, polygons, angles, and measurement. (I assume for purposes of this discussion that you are at least generally familiar with these concepts.) Furthermore, to compare the hexagonal grid to other possibilities the bees might use, in order to discern the reason they use hexagons (from an evolutionary point of view - I'm not suggesting the bees think about it), we'd need to be clear about some basic assumptions of *flat* plane geometry and the consequences of those assumptions.

One such assumption is that given any line l and any point P not on line l , there is a line through point P that is parallel to line l . This assumption reasonably captures part of what we mean by "flat." Observe that it is not true on a sphere. The shortest path between two points on a sphere runs along a *great circle* - a circle centered at the center of the sphere, making it as close to straight as a path on a sphere can be. Each great circle is the intersection of the sphere with a plane through its center. Notice as well that the sphere may be reflected onto itself in a great circle, just as a plane may be reflected onto itself in a line. Therefore, to a two-dimensional inhabitant of a sphere, the sides of the line are symmetrical: the line does not appear to be curving, because all of the curving occurs in the curvature of the sphere itself (and can only be seen from a three-dimensional perspective). One might say that a great circle curves *with* the sphere but not *in* the sphere. (By "in" we mean in the surface, not the space enclosed by the sphere in three dimensions.) Therefore, it is natural to interpret the "lines" on a sphere to be its great circles. Any two planes through the center of the sphere intersect in a line through the center of the sphere, and any such line intersects the sphere in a pair of *antipodal* points. (The North and South Poles give an example of a pair of antipodal points, if we visualize the earth as a sphere.) These antipodal points are on *both* great circles; thus, any two great circles intersect at two points. There are no parallel "lines" at all on a sphere! In fact, since lines in a plane can intersect in at most one point, whereas great circles intersect in two, great circles do not satisfy even the most fundamental assumptions that lines in a plane do!

The existence of a parallel line through any point not on a given line does not entirely

capture what we mean by “flat.” We further assume that if a pair of parallel lines is cut by a transversal, the alternate interior angles are congruent. (It is harder to imagine a surface that is not flat in this sense, and we won’t try to describe one here, but such surfaces do exist.) An important consequence of this assumption is that the sum of the degree measures of the angles of any triangle is 180° . From this result (theorem) we can further deduce that the corner angles of a regular triangle (usually called an equilateral triangle) must be 60° , the corner angles of a square (regular quadrilateral) must be 90° , those of a regular pentagon, 108° , and those of a regular hexagon, 120° . In general, we can prove the formula $\left(\frac{n-2}{n}\right)180^\circ$ for the corner angles of a regular n -gon. (Bear with me for a moment; we’ll see shortly how this relates to our question about the honeycomb.)

Remark. The word “equilateral” means that the sides are equal in length (that is, congruent). For polygons other than triangles, having congruent sides is not sufficient to make the polygon regular. The angles of an equilateral polygon with more than three sides need not be congruent. (An equilateral quadrilateral, for example, is called a rhombus. Draw some examples of rhombi that are not squares.) However, an equilateral triangle must be equiangular as well, a restriction that justifies our use of the word “equilateral” as equivalent to “regular” in this special case.

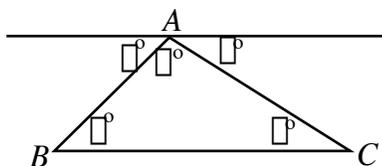
Example. Show that an equilateral triangle must be equiangular. On what assumption does this result rest? (By show we mean to prove informally. We cannot yet give a completely rigorous proof because we have not explicitly stated all of our assumptions and definitions.)

Solution: One of the basic assumptions of plane geometry is the side-angle-side (SAS) criterion for triangle congruence. (By definition, two triangles are congruent if there is a correspondence between their vertices for which the corresponding angles and sides are congruent.) Two sides of a triangle and the included angle determine the other side and angles; hence, given a correspondence between two triangles for which two pairs of corresponding sides and the included angles are congruent, the other corresponding parts must also be congruent. \square

From the SAS criterion, it immediately follows that the base angles of an isosceles triangle are congruent. (I expect you are familiar with this theorem.) For in $\triangle BAC$ (here we have labeled the vertices of an arbitrary triangle with capital letters to make discussion easier), if $AB \cong AC$, then $\triangle BAC \cong \triangle CAB$, by the SAS criterion. (That is, the triangle is congruent to itself with vertices A and B interchanged. We have used the symbol “ \cong ” as a shorthand to mean “congruent.”) Thus, $\angle A \cong \angle B$. Applying the base angle theorem twice to an equilateral triangle, we see that all of the angles are congruent.

Example. Consider a plane satisfying our second assumption: if two parallel lines are cut by a transversal, then the alternate interior angles are congruent. Show that the sum of the degree measures of any triangle in this plane must be 180° .

Solution. Given any triangle, call it $\triangle ABC$, consider a line through A that is parallel to line \overleftrightarrow{BC} , as in the figure below. Let the degree measures of angles $\angle A$, $\angle B$, and $\angle C$ be α , β , and γ , respectively. Since the alternate interior angles must be congruent, and the angle measure of a straight angle is 180° , we see that $\alpha + \beta + \gamma = 180$. \square



1.2 Exercises

1. Show that the sum of the interior (corner) angles of a polygon with n sides is $180(n - 2)^\circ$. (Hint: divide the polygon into triangles. An alternative method is to extend one end of each side and consider the sum of the angles supplementary to the interior angles; the total rotation as you travel around the polygon is 360° . It might not be so obvious that this fact also depends on the plane being flat. Can you see why?)
2. Use the result of the first exercise to prove that the angle measure of each corner angle of a regular n -gon is $\left(\frac{n-2}{n}\right)180^\circ$.

1.3 Back to the bees! Analyzing the honeycomb.

In order to analyze the honeycomb, we consider the abstract concept of a flat plane covered by a grid of congruent regular polygons. What shape can the polygons have? Since at least three corners must meet at any vertex of the grid, and the angles around a vertex must add up to 360° , we see that the only possible shapes are triangles, squares, or hexagons. The angles of a pentagon are too small to fit three around a vertex, but too large to fit four. The angles of any polygon with more than six sides are too large to fit three.

Now, compare the area to the perimeter for each of these figures. You will see that, for a given area, a hexagon has the smallest perimeter (followed by a square and then a triangle). Therefore, the bees can store the same amount of honey with much less beeswax in the walls of the honeycomb if they use hexagons! Making beeswax takes energy and material, and energy and material require food. By being more efficient, the bees are better able to survive and reproduce using the available food in their environment. It's not that bees can do geometric theory, but that natural selection has favored bees that make hexagonal honeycombs. Geometric theory has enabled us to see why!

1.4 Exercises

1. Calculate the area of an equilateral triangle as a function of its perimeter, p . (The units are not important; if some unit is chosen for length, the area will be in squares of that unit.)
2. Do the same for a square and regular hexagon.
3. Explain why, for a given area, the hexagon has the smallest perimeter, followed by the square and then the triangle.
4. Solitary insects and worms that build free-standing shelters for themselves typically build round cylinders. Why?

Remark. complete analysis would have to account for the forces on the walls, which would determine how thick they needed to be. These forces might be different for the differently shaped grids. But even our crude analysis gives insight into the pattern we see in a honeycomb.

1.5 Examining our assumptions.

We have seen that our assumptions are not true on a sphere. What if bees built spherical honeycombs instead of flat ones? Not surprisingly, the result about the angles of a triangle is not true on a sphere. The measures of the angles of a triangle must add up to *at least* 180° , and the sum can be much greater. Using either a model of a sphere or your imagination, you should readily be able to see that this is true. You should be able to construct an equilateral triangle whose angles all measure 90° , for example. (Hint: Imagining the sphere as planet Earth, start anywhere on the equator, travel directly to

the North Pole, turn 90° , return directly to the equator, and return along the equator to where you started.)

Because there is some flexibility in the angle size of a regular polygon on a sphere, we can cover a sphere with a grid of triangles, squares, or pentagons. However, hexagons cannot be used, nor can polygons with more than six sides, because the angles are too large. It can be shown that the angles of a polygon determine its size. (The larger the polygon, the larger the angles, as you can convince yourself by looking at some examples.) Therefore, there is only one way to make each possible type of grid on a given sphere. There must be either four triangles, with three meeting at each vertex, eight smaller triangles, with four meeting at each vertex, 20 even smaller triangles, with five meeting at each vertex, six squares, with three meeting at each vertex, or 12 pentagons, with three meeting at each vertex. These grids give rise to the five Platonic solids: the tetrahedron, octahedron, icosahedron, cube, and dodecahedron, respectively. These are the only possible regular polyhedra. (Why can't four squares or pentagons meet at a vertex?)

1.6 Extra credit exercises!

1. Is the theorem about the base angles of an isosceles triangle true on a sphere? Explain your answer. (Hint: Is the SAS criterion true on a sphere? What makes the SAS criterion work? Think about the motions that move the points of a plane or sphere around without altering the distances between them.)
2. Is the SAS criterion true on a surface that is rounded in some places and flat in others?
3. Which grid would the bees use if they built spherical honeycombs? Why?
4. Why do you think bees do not make spherical honeycombs?

1.7 Mathematics is abstract by nature.

To analyze the honeycomb, we made an abstract model of it that captured its most important properties, a sort of “ideal” honeycomb. Good as the bees are at building – and they are very good – the cells of a real honeycomb are not perfectly regular, nor do they have exactly the same volume. Furthermore the walls of the honeycomb are not of exactly uniform thickness. (This variation is not all due to error. The interior walls

of the cells are intentionally rounded; therefore, the cell walls are thicker at the corners. The bees also build cells for raising young, which have different sizes depending on the type of bee being raised.) Our assumption that the cells are hexagonal is more closely realistic; my entomologist colleagues did not know of any instances in which the bees used a different number of sides for any cell. It is not surprising that this assumption accords most closely with reality, since it is fairly stable with respect to variation in the sides and angles: if the bee turns somewhat more or less than 120° at the corners as it proceeds around, or if it makes some sides a little shorter or longer, the shape will still close up into a hexagon.

When we create an abstraction from a collection of physical examples, we *assume* that some properties hold exactly, even if they are only approximately true in the examples. In addition, we ignore many details, which may vary. When we abstract from mathematical examples, we also ignore many variables while focusing on essential commonalities. For example, a collection of triangles may have various shapes and sizes, but they all have three vertices and three sides. It is these two simple properties that define the concept of triangle, of which there is a great variety of examples.

Examining our assumptions leads to new possibilities. What if we made different ones? What if conditions were different?

Even the simplest mathematical concepts are inherently abstract. Consider, for example, a natural number such as three. We know what three oranges are, or three cats, or three dogs, but what is the quality of “having three elements” that these collections of objects have in common? A set has three elements if it can be put in one-to-one correspondence with the set $\{1, 2, 3\}$. (This process is, after all, how we count.) However, this observation doesn’t help much, since the elements of this set are numbers, which we are trying to define! What are these numbers? The elements 1, 2, and 3 are not just symbols, because if we wrote these numbers in Japanese they would still be the same numbers and serve the same purpose. It took mathematicians a long time to come up with a satisfactory, universal answer to this question, which we introduce toward the end of this course.

We cannot completely understand the properties of an abstract concept by experiment. For one thing, we cannot test all of the infinitely many possibilities. Even if we could, our results would give us no insight into the relationship between the properties we discovered and the defining properties we assumed, and it is exactly these relationships that constitute explanation and understanding. Therefore, we must logically *deduce* the consequences of our assumptions. It is often surprising, and enlightening, how much can be deduced from very simple assumptions. For example, the two simple properties that define a triangle, along with the SAS criterion and some simple assumptions about lines and points, have as consequences that the bisectors of the angles of any triangle

are concurrent (that is, they all intersect at a single point). Since we do not need any assumptions that depend on our surface being flat, this result is true on a sphere as well.

In order to make deductions correctly, we must do two things:

- Precisely specify our assumptions and definitions; and
- Establish effective steps that determine whether or not a conclusion has to be true under the assumed conditions.

Learning to do these two things is a central purpose of this course! Another central purpose is to learn to communicate our reasoning clearly to others and to read mathematical texts effectively, extracting their essential ideas and the relationship of those ideas to what we already know. Along with deductive reasoning, metaphor and analogy play a role in communicating mathematical ideas and nurturing an understanding of the mathematical landscape that illuminates formal derivations (which anyone, even the most devoted “geek,” would find insufferably boring by themselves!).

Please note that truth in mathematics is always conditional. When we say a theorem is “true,” we mean that it is true on condition that the assumptions of the theory under discussion are true. We make these assumptions explicit initially, but of course it would not be practical to include them as part of the statement of every theorem; instead, they may be thought of as a sort of definition of the whole system under study. We will not deal with notions of “universal” truth (except in the context of a particular specified “universe”). Such notions belong to religion or philosophy (if they belong anywhere at all), not mathematics.

The concept of validity is also used in mathematics, and is related to the idea of truth, but with slightly different connotations. An argument is *valid* if it is logically correct; a conclusion is valid if it has been derived from the stated assumptions by a valid argument.

Some mathematicians prefer to avoid using the word “true,” in order to emphasize the conditional nature of mathematics. Indeed, what would it mean for a mathematical statement, which refers to purely abstract products of our collective imagination, to be true? In the final analysis, mathematics is a purely formal system that is only given life and meaning by the way it models aspects of our experience. However, I find it awkward and unnatural to completely avoid the word “true.” It is often convenient to refer to statements as “true” for purposes of discussion, as long as we keep in mind the contextual nature of this “truth.”

1.8 The dance of abstraction and concreteness.

Making mathematical models is a crucial way of understanding the world. Without them, we would have only collections of facts and observations without any understanding of the relationships among them. Without them, we cannot understand how or why things work the way they do, or make any but the most simplistic predictions.

This is not to say that the only purpose of mathematics is its application to science and other disciplines. Many people find the ideas of mathematics to be beautiful and inspiring in themselves. It gives mathematicians joy to pursue mathematics for its own sake, regardless of any application their work might have.

It is also worth noting that often work that was initially thought of as entirely “pure” does turn out to have valuable applications, such as the use of number theory to develop codes for the secure transmission of information. These codes are used everyday, whenever we make a financial transaction over the internet.

By deepening our fundamental understanding of patterns, structure, and logic, we make possible applications that we could not have even conceived, much less achieved, before. Practical application is sort of like happiness: you achieve it best if you do not pursue it too directly!

As mathematics develops, there is a drive for increasing theoretical abstraction, which allows for more general results. This drive operates in tension with the competing drive to find and explore concrete examples that inform our intuition (sometimes by reflecting our real world experience, sometimes by referring to abstractions with which we have become familiar and comfortable), giving relevance and meaning to the concepts we have developed. This tension is never-ending and inevitable.

A closely related dynamic tension exists in mathematics between the *constructive* and *synthetic* approaches to its development. The constructive approach strives to represent mathematical objects as particular sets, which gives them a concrete character, although not exactly in a physical way. The idea is to “construct things from scratch,” so to speak, deriving their essential properties from the constructed representations. A classic example is the representation of the numbers $0, 1, 2, 3, \dots$ by the sets $\emptyset, \{\emptyset\}, \{\emptyset, \{\emptyset\}\}, \{\emptyset, \{\emptyset\}, \{\emptyset, \{\emptyset\}\}\}, \dots$, where the symbol “ \emptyset ” denotes the unique set containing no elements, known as the *empty set*. This representation has two important and noteworthy properties: first, the set representing each number is constructed in a systematic way from the preceding one, by simply including it along with all of its elements (the previous representatives), so there is a natural means of continuing the sequence; second, each representative has the number of elements that it represents. In contrast, the synthetic approach proceeds by simply assuming the existence of concepts satisfying certain properties and exploring what must be true about them as a consequence of

these properties. For example, we begin our study of algebra by assuming the existence of a system of real numbers with specified operations and properties. A classic example of the synthetic approach occurs in geometry, in which we begin by assuming there are things called points and lines, which we cannot define, and some relationships among them, which we also cannot define, such as that a line may *pass through* a point, which would then *lie on* that line, or that one point may lie *between* two other points on a line. (We do visualize points and lines, perhaps thinking of them as dots and lines on paper, but these representations are physical, not mathematical, and only approximate.) We assume additionally that these relationships satisfy certain properties, such as that given two distinct points, there is exactly one line passing through both of them. (Otherwise, the theory would be a meaningless free-for-all!) The properties we choose are abstracted from our visualizations. From our undefined concepts, we can define new ones, such as that of a segment: given two points, A and B , *segment AB* is the set of points containing A , B , and all the points between A and B . We can in turn assume the existence of relationships between these defined concepts, such as congruence, explore their properties, and so forth. The rest, as they say, is history!

Mathematics moves fluidly between construction and synthesis, starting at various points and moving both towards greater complexity and towards greater fundamental understanding and unification. New concepts are constructed from synthetic ones, as in the example of segments, constructed from the undefined concepts of *point* and *betweenness*, and they in turn may motivate new synthetic concepts, such as congruence. After exploring a synthetic system of axioms, we may seek a foundation of more basic concepts from which a model satisfying these axioms can be constructed. Alternatively, after exploring a construction, we may seek a set of axioms that underly and explain its properties, which might in turn apply to other constructions. To use geometry as an example, we might initially take the synthetic approach pioneered by Euclid and brought to culmination by Hilbert,¹ but eventually construct mathematical representations of points and lines using a coordinate system. Among other things, we might then study the transformations of the coordinate plane that take segments and angles to congruent ones. The common properties exhibited by these collections of geometric transformations might inspire us to consider some of the structures of abstract algebra, such as groups, defined merely as sets with certain operations, satisfying certain properties. (Now, rather than abstracting from the physical world, we are abstracting from mathematics, focusing on some properties while ignoring other, more specific ones.) The foundational concepts

¹David Hilbert, probably the most esteemed mathematician of the late nineteenth and early twentieth centuries, developed a complete set of axioms for plane geometry in the spirit of Euclid, filling in the implicit assumptions that Euclid had made but left unstated.

of mathematics, codified in the axioms of set theory and logic, are synthetic: you have to start somewhere.

For some sets of properties, there are many examples satisfying them, but if we assume enough properties, essentially one representation will remain, in that any other will have exactly the same structure, just described in different terms. Once you assume enough, there is no room for variation. For example, there is essentially one representation satisfying all of the properties of Euclidean geometry, which we call the Euclidean plane. You are familiar with its standard construction from high school algebra: the points are pairs of real numbers, and the lines are the solution sets of linear equations in two variables.

1.9 A note on how this course will proceed.

Courses on mathematical methodology and structure tend to falter in one of two ways. At one extreme, one can avoid all ambiguity by adopting a high level of formality from the start, insisting for example that we must define the Cartesian product of two sets (don't worry if you don't know what that is yet) before we can talk about familiar binary operations such as addition. The result is tedium, along with the loss of any appreciation for how or why anyone came up with these formal notions in the first place. Proceeding in this fashion is akin to trying to learn grammar before you learn to talk (although once you can talk, you most certainly should refine your understanding of language by studying grammar). I am sure you have no trouble understanding, without reference to Cartesian products, that addition is called an operation because it takes a pair of numbers and associates to them a third number, their sum, or that it is called binary because it operates on two numbers at a time. Yet, at the opposite extreme, one can begin quite imprecisely, assuming that a vague familiarity with the properties of numbers is sufficient to start "proving" things. The result is confusion, as it is not clear to students what may be assumed (that is, used in a proof) and what may not.

We will follow neither of these approaches, nor even a compromise between them. Instead we will begin with familiar concepts, described at a comfortable level of abstraction and formality, but with precision. As we proceed, we will examine these concepts in increasing depth and also see them within an increasingly general context. We will temporarily have to spell out some of our assumptions in non-technical language. (Eventually we will see that they all boil down to the existence of certain sets; however, whereas this technical reduction provides clarity, insight, and a certain kind of satisfaction, it is not necessary for unambiguous communication.) We may make some assumptions provisionally, as mathematicians often do when interested in exploring their consequences,

informally justifying or motivating these assumptions and proving them later based on more fundamental ones. The important point is that we will attempt to be both explicit and exact about our current assumptions at all times. These and only these explicit assumptions, along with previously proven results, may be used in proofs.

We begin with some motivation, in which we ask you to examine familiar procedures to determine what unstated, implicit assumptions and processes of reasoning, so generally accepted we may not even notice them, are required to justify what we do. We hope that this exercise will help you understand and appreciate the mathematical theory it is used to motivate.

Chapter 2

Motivation: Area Revisited

In which we examine the fundamental principles of measure and derive some formulas from them.

2.1 What is Area?

Continuing in a geometric vein, let us consider area on a flat surface. Rather than begin with formulas by which it may be calculated, let us instead consider what we mean by area: what area is and what it should do. From these assumptions, we will then *deduce* formulas for calculating it in some instances. In this manner, we will better understand the *meaning* behind the formulas. This understanding is a crucial aspect of mathematical thinking, one with which you may not have as much experience as you should!

Area is a positive numerical quantity associated to a region of the surface. Thus, area is a type of *measure*; the number associated to a region is its measurement. What properties should such a measure have?

The numbers associated to two regions should allow us to compare the sizes of the regions, but we must be precise about what we mean by “size,” since we can measure and compare sizes in different ways. If we wanted to know how long it would take to walk around the boundary of a region, we would measure its size by the distance we would have to travel, but that is not what we mean by area; the distance around the boundary of a region is its *perimeter*. If we want to know how much wheat we could grow in the region, the area measurement (along with the yield for wheat) should tell us. A long, skinny region could have the same area as a shorter, fatter region, but a much larger perimeter.

As with length, the area of a region should depend only on its size and shape: con-

gruent figures should have the same area. Mathematicians say that area is a *geometric invariant*, meaning that it does not vary among figures that are geometrically the same, in the sense of being congruent. Perimeter is also a geometric invariant.

Obviously, the property of being a geometric invariant is not enough in itself to determine area, since it is true of perimeter as well. For a second property, we can also take a clue from the properties of length. Note that if a segment (of either a straight line or a curve) is divided into parts that do not overlap except at their endpoints, its length is the sum of the lengths of its parts. Similarly, if a region is divided into parts which do not overlap except along their boundaries, its area should be the sum of the areas of the parts. If I can grow 10 bushels of wheat in a field and 7 bushels in a neighboring field, then I can grow 17 bushels in the region covered by both fields. (Notice that area differs from perimeter in this respect; perimeter does not have this property, because a shared edge would count only once when the regions are combined.)

We are really talking about analogous measures in different dimensions. Length is one-dimensional, so zero-dimensional objects, namely individual points, have no length. If the parts of a segment overlap only at their endpoints, the length of the whole segment is the sum of the lengths of its parts because no length is shared. Area is two-dimensional, so zero- and one-dimensional objects, namely points and segments, have no area. If the parts of a region overlap only in one-dimensional objects, then its area is the sum of the areas of its parts because no area is shared.

Remark. In higher mathematics, it turns out to be better to consider *oriented* measures. Oriented length is positive in one direction, negative in the other, and this notion can be extended to higher dimensions. But that is another story.

The two properties discussed above are fundamental, but alone they are not sufficient to determine a unique area measure. This is because any measure with these properties is *scalable*: if we double, triple, halve, or scale all of the measurements by any other factor, the two properties clearly still hold, and the comparison of regions remains valid. To determine a particular area measure, we must choose a unit by convention. (This is also true for perimeter, of course; you get different measurements in meters and feet, for example.)

Our area unit will depend on the unit chosen for the more basic measure of distance. Once this unit is chosen, say centimeters (cm), we define the area of a rectangle to be the product of its length and width: if R is a rectangle with length l cm and width w cm, then its area is $A_R = lw$ cm². “Wait,” you say! “You just threw a formula at us; I thought we were going to deduce the formulas!” Well, yes, we are going to deduce all of the other formulas, but in order to establish an area measure we need to assume one formula.

Our assumed formula must be consistent with our first two assumptions. Clearly, it is. Congruent rectangles have the same length and width, hence our formula gives them the same area. And if a rectangle is divided into two smaller rectangles by a segment parallel to a pair of its sides, the area given by our formula for the whole rectangle is the sum of the areas given by our formula for the two smaller rectangles. (See Figure 2.1.) This fact depends on the distributive property of addition and multiplication, which one might view as being assumed in order to make number operations consistent with our intuitive understanding of area! By seeing that our formula is a natural choice having the two desired properties, we understand the meaning behind it!

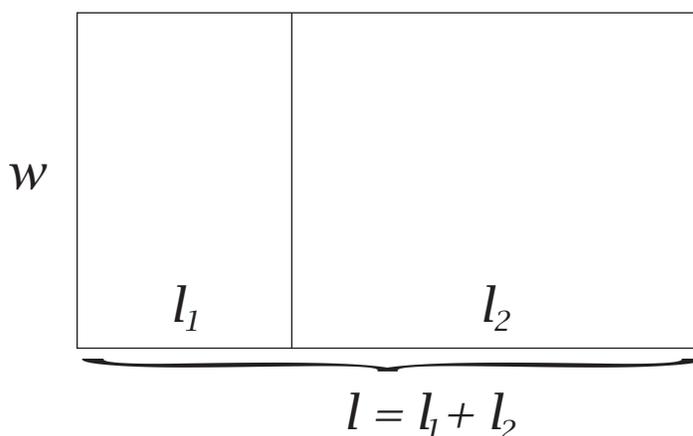


Figure 2.1: $lw = (l_1 + l_2)w = l_1w + l_2w$

Remark. This method of establishing a unit does not work on a curved surface, such as a sphere, because rectangles do not exist on a curved surface. We can certainly measure area on a curved surface, but we need a different method of establishing a unit.

From the area formula for a rectangle, we can easily deduce the area formula for a parallelogram. Observe that by dropping a perpendicular from an obtuse corner to the base, we divide a parallelogram of base b cm and height h cm into a right triangle and a trapezoid; hence, its area is the sum of the areas of these two figures. Observe further that, since the height of the parallelogram is measured perpendicular to the base, a rectangle with sides of length b cm and h cm may also be divided into a triangle and trapezoid that are congruent, respectively, to these; hence, the rectangle's area is given

by the same sum. Since the area of the rectangle is known to be bh cm², the area of the parallelogram must also be bh cm². (See Figure 2.2.)

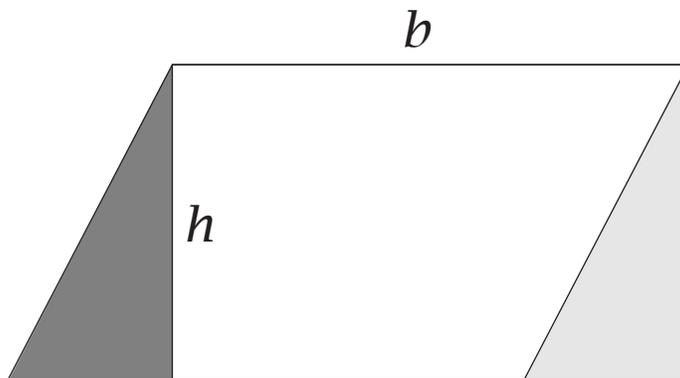


Figure 2.2: The area of a parallelogram of base b and height h is bh .

2.2 Exercises

1. Deduce that the area of a triangle of base b units and height h units is $\frac{1}{2}bh$ squared units. (Hint: Construct a parallelogram from two copies of the triangle, one of which is rotated by 180° from the other; since the triangles are congruent, each has half the area of the parallelogram.)
2. Deduce that the area of a trapezoid with parallel bases of length b_1 units and b_2 units and height h units is the height times the average of the bases: $\frac{1}{2}(b_1 + b_2)h$ squared units. (Hint: Divide the trapezoid into two triangles; alternatively, construct a parallelogram from two copies of the trapezoid; alternatively, divide the trapezoid into a parallelogram and a triangle; alternatively, divide the trapezoid into a rectangle and two right triangles. There is more than one way to skin a cat!)
3. If the length of a rectangle is 3 meters and its height is 51 centimeters, could we sensibly describe its area as 153 meter-centimeters? If so, what region naturally represents one meter-centimeter?

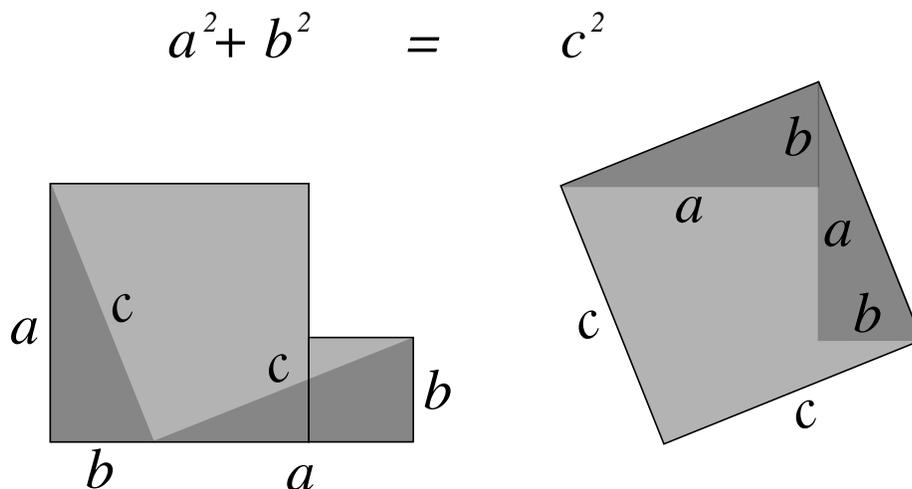


Figure 2.3: The Pythagorean Theorem

4. The Pythagorean Theorem states that if a right triangle has legs of length a units and b units and hypotenuse of length c units, then $a^2 + b^2 = c^2$. Referring to Figure 2.3, deduce the Pythagorean Theorem from the properties of area.
5. Like most theorems, the Pythagorean Theorem is *conditional*: it has a *hypothesis*, which imposes conditions under which its *conclusion* must hold. As the theorem is stated above, the conditions are all mixed together, conveyed via the words “right,” “legs,” and “hypotenuse,” along with the manner in which the sides are labeled. The the geometrically significant condition is clearly that the triangle has a right angle; while it is certainly important to label the sides correctly in the formula, this consideration is merely technical.

It is often the case that we can clarify the meaning of a theorem by laying out the technical conditions in a preliminary exposition in order to isolate the more meaningful hypotheses in the final statement. Let us restate the Pythagorean Theorem in this manner: Given any triangle, let its vertices be labeled A , B , and C , and let the sides opposite A , B , and C have lengths a , b , and c units, respectively. If $\angle C$ is a right angle, then $a^2 + b^2 = c^2$.

The *converse* of a conditional statement is obtained by exchanging its hypothesis

and conclusion. Keeping the same preliminary conditions, a meaningful converse to the Pythagorean Theorem is immediate from the version above: If $a^2 + b^2 = c^2$, then $\angle C$ is a right angle.

Is the converse of the Pythagorean Theorem true? If so, outline how you would prove it, stating any additional basic geometric results (such as congruence criteria) you would need. If it is false, give a counterexample. (A *counterexample* is a specific example for which the hypothesis is true but the conclusion is not, thereby proving that the statement in question is false.)

6. There is also an unstated condition preliminary to the Pythagorean Theorem as we deduced it: the triangle must lie in a flat surface. Without this condition, the argument using our area formulas would not apply.

Just because a particular argument applies only under certain conditions does not mean a statement is not also true under other conditions; it might simply be that a different argument is required. (Perhaps a more general argument that applies under both sets of conditions can even be found.)

In this light, is the Pythagorean Theorem true on a sphere? If so, outline how you would prove it, stating any properties of spherical geometry you would need. If it is false, give a counterexample. (Hint: consider the discussion of spherical triangles in Section 1.5.)

Remark. Many treatments of area take triangles rather than rectangles as the fundamental regions, since any polygon can be divided into triangles. This approach is necessary on the sphere, where rectangles don't exist.

Remark. Calculating the area of circles and other regions with curved boundaries requires the concept of limit, as defined and studied in calculus. This approach goes back at least to Archimedes in ancient Greece, but its precise logical development occurred much later. Calculus is also needed to calculate the areas of regions on surfaces of non-uniform curvature.

2.3 What have we left out?

I hope the arguments above are clear and convincing; nonetheless, much has been left out. You probably didn't notice because the material is so familiar. Can you see, for example, how we used the fact that the surface under consideration is flat? (Hint: It has to do with parallel lines and angles. As mentioned previously, rectangles, defined as quadrilaterals with four right angles, don't even exist on a sphere.)

Not only are there unstated details in the arguments, there are unstated assumptions whose proof requires theoretical development we have not covered. In common parlance, there are logical *gaps* (sort of leaps of faith, except that perhaps you didn't notice you were leaping). For example, the argument for a parallelogram requires knowing that the opposite sides are congruent. This fact is not part of the definition, which only requires the opposite sides to be parallel. That they must also be congruent can be deduced by dividing the parallelogram into two triangles; however, these triangles share only one side, and we know nothing about the other pairs of sides. (Remember, we are trying to deduce logically that they are congruent; we cannot assume they are.) Thus, to prove the two triangles are congruent, we need another criterion, the angle-side-angle (ASA) criterion. (We know the corresponding angles are congruent because they are alternate interior angles for parallel lines cut by a transversal, namely the diagonal that divides the parallelogram into triangles.) Do we need to make ASA a separate assumption, or can it be deduced from more fundamental assumptions, such as SAS? How do we know this new criterion is even consistent with our other assumptions? In fact, ASA can be deduced from SAS and more fundamental assumptions. Although doing so is not particularly difficult, it is somewhat involved and beyond the scope of this course.

As another example, consider the perpendicular dropped from the obtuse corner of the parallelogram. How do we know we can drop a perpendicular from any point to any line, and how do we know that in this case it lands on the base, rather than outside the parallelogram? Proving this last point requires a major result of plane geometry called the Exterior Angle Theorem. This is not a trivial point, since the Exterior Angle Theorem is not true on a sphere. (Fortunately, there are no parallelograms on a sphere, either!) Again, the theoretical development required is not particularly difficult, but it takes considerable time and is beyond the scope of this course. You will have the opportunity to fill these gaps by taking a course in classical geometry.

Logical gaps are intrinsic to the process of learning and discovery. We must often make conjectural leaps in order to see the framework of a theory, either coming back to prove these conjectures later or trusting the work of others. If they turn out to be false, we start again with a new approach, informed by the insights gained from seeing our errors. Indeed, if we tried to proceed step by logical step without ever trying to see further ahead, we would never get very far. It would be like walking while looking down at one's feet!

It is neither practical nor necessary to fill in every logical step as we proceed; however, it is very important to recognize when we are leaping over gaps and consider the logical steps that might bridge them.

Chapter 3

Motivation: Solution Sets

In which we use the example of solving an equation to introduce some fundamental concepts, language, and notation used in set theory and logic.

You are familiar with solving equations. Let us look at some simple equations in order to understand more precisely what it means to solve one and by what process we do it.

3.1 What does solving mean exactly?

Consider the following simple equation:

$$x + 2 = 5.$$

An equation such as this one is a statement in mathematical language. Translated into English, this statement says, “When two is added to some number, let’s call it x , the result is five.”

Because the equation involves an unknown quantity, the number represented by the variable x , it implicitly asks a question: What number could x be? When x is replaced by a specific number, such as 3 or 4, the resulting equation must be either true or false (but not both). As you well know, the equation $3 + 2 = 5$ is true, whereas $4 + 2 = 5$ is false. (As emphasized in Chapter 1, this notion of truth or falsity is conditional, based on the assumptions we make about numbers and the operation of addition. Indeed, it also requires the *definitions* of 2, 3, 4, and 5 in terms of 1: $2 = 1 + 1$, $3 = 2 + 1$, $4 = 3 + 1$, $5 = 4 + 1$. By definition of 4, $4 + 1 = (3 + 1) + 1 = 3 + (1 + 1)$. By definition of 2, $3 + (1 + 1) = 3 + 2$, so $3 + 2 = 4 + 1 = 5$, where the final equation follows from the definition

of 5. Do you remember the name for the property we used to regroup the numbers along the way? We have not yet made our assumptions about numbers explicit, but they are based on our innate intuition about what whole numbers and addition mean. This intuition is so fundamental that it is not unique to humans; other animals, such as crows and dolphins, can also count and add small whole numbers.)

In the study of logic, a statement that must be either true or false is called a *proposition*; formally, we say a proposition has a unique *truth value*. Not all statements are propositions. For example, “Mozart wrote beautiful music” is a statement, but we cannot assign it a unique truth value; whether or not it is true is a matter of individual opinion.

Because an equation has a unique truth value for any number that is substituted for the variable, it defines a *set*, meaning a collection of things that can be identified, called its *solution set*: the numbers for which the equation is true are in the solution set, the others are not. Solving the equation means finding the most transparent and direct description of this set: instead of merely knowing that it is the set of numbers, x , for which $x + 2 = 5$, we prefer to also know that this set contains exactly the number 3 and no other numbers.

The things contained in a given set are called its *elements*. Symbolically, if S is a set, we write $x \in S$ to denote that x is an element of (or “belongs to”) S . By the nature of the concept, a set is completely characterized by its elements; thus, if two sets have exactly the same elements, they are equal (that is, they are the same set). An equation defines a set in terms of a proposition about its elements. More symbolically, we write:

$$\{x : x + 2 = 5\}.$$

The brackets tell us we are expressing a set, and the colon translates into English as “such that.” Thus, the expression above translates into English (or, more accurately, “Mathlish”) as “the set of all x such that $x + 2 = 5$.” Many authors use the symbol “|” instead of the colon; either symbol is equally acceptable.

Our solution defines the same set as an explicit list of elements (in this case, a list of length one):

$$\{x : x + 2 = 5\} = \{x : x = 3\} = \{3\}.$$

A proposition about a variable is called an *open* proposition. Two open propositions are considered to be *equivalent* if they are true for exactly the same set of things, that is, if they have the same solution set. Thus, the open propositions $x + 2 = 5$ and $x = 3$ are equivalent. The sets $\{x : x + 2 = 5\}$ and $\{x : x = 3\}$ that are described by these equivalent open propositions are equal because they contain the same elements, just

identified in different ways. In summary, equivalent open propositions describe equal sets.

If we can find an open proposition equivalent to the original equation that explicitly lists the elements of the solution set, that is, a proposition of the form “ $x =$ this number or $x =$ this number . . . or $x =$ this number,” then we have solved it. Note that the order of the list of elements for a set makes no difference; it doesn’t even matter if an element is repeated. We are only concerned with what is in the set. Note in addition that the list might be empty, as for the equation $x + 2 = x$. In this case the solution is the *empty set* - the unique set with no elements at all, denoted by the symbol \emptyset . Although the solution to the simple equation $x + 2 = 5$, which is only one step away from a direct description, is obvious, many equations are more difficult to solve; consequently, knowing their solutions is commensurately more helpful.

A theme of this book is that mathematics, although special, is not that different from the rest of life, and this situation is no exception. For example, the proposition, “Some person visited my friend John yesterday evening,” invites an attempt to figure out, by logical investigation, who it was. If we find that this proposition is equivalent to, “Some person is my friend Janet or some person is my friend Todd,” we have solved the mystery (and can move on to wondering what Janet, Todd, and John might have talked about).

Before examining what is involved in finding an explicit list of solutions to an equation, we should pause and be a little more careful about just which numbers we are considering as possible values of x . In fact, since we haven’t placed any restriction on what x could be in our description of the set $\{x : x + 2 = 5\}$, x could be any mathematical object at all (whatever all the possible mathematical objects might be - another good and serious question); although the defining equation only makes sense for objects to which 2 could be added, this sloppiness is clearly unsatisfactory and we will henceforth studiously avoid it!

I assume you are sufficiently familiar with the systems of natural numbers, integers, rational numbers, real numbers, and complex numbers to do at least simple computations; we will discuss all of these number systems with great care in future chapters. (By *system* we mean a set of numbers together with the operations, such as addition or multiplication, that apply to them.) Let us review, since we need to use them, the standard symbols that represent the sets of numbers in these systems: the set of natural numbers is denoted by \mathbb{N} ; the set of integers is denoted by \mathbb{Z} , the set of rational numbers is denoted by \mathbb{Q} , the set of real numbers is denoted by \mathbb{R} , and the set of complex numbers is denoted by \mathbb{C} . Each of these systems expands on the previous one. Authors vary as to whether or not zero is included in the set of natural numbers. We will include it: $\mathbb{N} = \{0, 1, 2, 3, \dots\}$.

For the simple equation $x + 2 = 5$, it does not make any difference in which system

we are working, since we get exactly the same solution in all of them; however, for many equations the system does make a difference. For example, the equation $2x = 5$ has no solution in the system of natural numbers, although its solution set in each of the other systems is $\{\frac{5}{2}\}$. The equation $x^2 = 2$ has no rational solutions, but its solution set in either the real or complex number system is $\{-\sqrt{2}, \sqrt{2}\}$. (We will prove later that $\sqrt{2}$ is not rational; that is, it cannot be written as a ratio of two integers. This fact was discovered in ancient Greece by the school of Pythagoras.) In the natural, integer, rational, and real number systems, the equation $x^3 = 1$ has a unique solution, 1, but in the complex number system it has three: 1, $\frac{-1-i\sqrt{3}}{2}$, and $\frac{-1+i\sqrt{3}}{2}$.

We can specify the set of possible values for x in one of two ways. One way is to include it as an additional requirement in the proposition that defines the set, as in:

- $\{x : x \in \mathbb{R} \text{ and } x + 2 = 5\} = \{3\}$
- $\{x : x \in \mathbb{R} \text{ and } x^3 = 1\} = \{1\}$
- $\{x : x \in \mathbb{C} \text{ and } x^3 = 1\} = \{1, \frac{-1-i\sqrt{3}}{2}, \frac{-1+i\sqrt{3}}{2}\}$

This approach, while logically correct, may seem a bit awkward because the two requirements in our description, $x \in \mathbb{R}$ and $x + 2 = 5$, play rather different roles; the former is far more general than, and provides a context for, the latter. As we learned in our discussion of the Pythagorean Theorem, it often improves clarity to get contextual assumptions out of the way up front, which we can do with the following alternative, and slightly briefer, construction:

- $\{x \in \mathbb{R} : x + 2 = 5\} = \{3\}$
- $\{x \in \mathbb{R} : x^3 = 1\} = \{1\}$
- $\{x \in \mathbb{C} : x^3 = 1\} = \{1, \frac{-1-i\sqrt{3}}{2}, \frac{-1+i\sqrt{3}}{2}\}$

Note that the symbol “ \in ” may be used as either a verb, as it is when it appears in a defining proposition (“such that x belongs to \mathbb{R} ”) or a preposition, as it is when it is used to restrict the values of x to a particular set (“ x in \mathbb{R} such that ”). Mathematical usage liberally incorporates grammatical flexibility of this kind.

Now that we have thoroughly discussed both the context and meaning of solving an equation and introduced the language needed to formally describe solutions, let us turn to the process of finding the solution.

Remark. Ordinary language is often used in preference to formal and symbolic mathematical descriptions. We more often write, “The real solution to the equation $x + 2 = 5$ is 3,” than “ $\{x \in \mathbb{R} : x + 2 = 5\} = \{3\}$.” In many situations, formal mathematical language is just not as natural, and it is limited in its expressive power to convey the full subtlety of ideas: the analogies, comparisons, and connections that occur behind the scenes in our thinking. (But notice that we do prefer the symbolic language of algebra, which is so clear and concise, for expressing numerical operations and relationships.) In mathematical formalism, there are no subtexts implied by a particular choice of words! Its virtues include uniformity and precision, which can help in logically tricky situations.

3.2 The process of solving an equation

For the simple equation above, we can summarize our process simply as subtracting 2 from both sides; however, this pat description hides several properties of numbers and number operations that are being used. Aside from knowing how to add and subtract explicit numbers correctly, not to mention the definitions of number names alluded to earlier, we also needed the following facts:

- The operative quality of addition: the sum of two numbers is determined by the numbers that are added. If we start with two equal numbers (albeit expressed in different ways) and add the same number to each, we will get the same sum in both cases.
- The definition of subtraction in terms of additive inverses: subtracting 2 is the same as adding -2 . (Since -2 is not a natural number, we technically need to define subtraction differently in the natural number system, and of course it is not defined for all ordered pairs of inputs; it is easiest just to work in a larger system that has additive inverses when solving equations.)
- The associative property of addition: $(x + 2) + (-2) = x + [2 + (-2)] = x + 0$. We needed first to describe subtraction in terms of addition because the operation of subtraction is *not* associative; for example, $(5 - 2) - 3 \neq 5 - (2 - 3)$. The associative property is relevant because addition is a binary operation: one can only add two numbers at a time.
- The defining property of 0 as the additive identity: $x + 0 = x$.

Let us examine in detail the logical structure, meaning, and significance of each of these assertions.

The operative quality of addition is a *universal* property: it applies to all numbers in the system. Taking the real numbers as our system, for example, the precise statement of this property is that, for all real numbers x, y, x' , and y' , if $x = x'$ and $y = y'$, then $x + x' = y + y'$.

Words such as “and,” “or,” “not,” and “implies,” and the pair “if ... then ...” act on propositions (open or not) to create new propositions whose truth value in any instance follows from that of their components. For example, a proposition composed of two propositions joined by “and” is true exactly for the instances in which both of its components are true. The new proposition that results is called the *conjunction* of its two components. As another example, if the phrase “it is not the case that” is inserted before a proposition, the resulting proposition is true exactly for the instances in which the original proposition is false. The resulting proposition is called the *negation* of the original. Since they act on propositions, conjunction and negation may be thought of as logical operations; conjunction is binary, acting on a pair of propositions, whereas negation is unitary, acting on a single proposition. (Note the similarity in form between such phrases as “the sum of x and y ” and “the conjunction of $x = x'$ and $y = y'$.” Similarly, note the similarity in form between “the cubed root of x ” and “the negation of $x = x'$.” Addition is a binary operation on numbers; taking the cubed root is a unitary operation.) There are two other logical operations, both of them binary; these are, of course, the operation represented by “or” and the one represented by “implies” (which is also represented by “if ... then ...”). A proposition that simply asserts a relationship between two things without involving any logical operations is called *atomic*. The negation of an atomic proposition can generally be abbreviated with a slash: for example, the negation of the proposition that $x + 2 = 5$ is the proposition that $x + 2 \neq 5$.

A statement of the form “ $x > 2$ or $x^2 > 4$ ” is called the *disjunction* of its two components. The mathematical convention is that disjunctions are never exclusive; that is, a disjunction is true if either or both of its component propositions are true, false only if both component propositions are false. For example, the proposition that $x > 2$ or $x^2 > 4$ is true for $x = 3$ as well as for $x = -3$. (If a flight attendant asks if you want a blanket or a pillow, he is using disjunction in the non-exclusive sense, since of course you can have both - and probably want both if you plan to go to sleep. If he asks if you would like coffee or tea, the sense of the disjunction is exclusive, since the rules allow only one drink at a time.) Note that it is impossible to find an instance in which $x > 2$ is true but $x^2 > 4$ is false, a consideration that will be of interest shortly.

The proposition that if $x = x'$ and $y = y'$, then $x + x' = y + y'$, is called *conditional* because it sets up a condition, in this case that $x = x'$ and $y = y'$, under which a consequence, in this case that $x + x' = y + y'$, is held to be true. Conditional propositions, also called *implications* are the most important in logic, because they are the way “one

thing leads to another.” For example, as we recently noted, the condition that $x > 2$ ensures in any instance that $x^2 > 4$, so the proposition that $x > 2$ implies $x^2 > 4$ is true for all real numbers x . There are many ways to phrase a conditional proposition; all of the following represent the same logical combination of the atomic propositions $x > 2$ and $x^2 > 4$:

- If $x = 2$, then $x^2 = 4$.
- $x = 2$ only if $x^2 = 4$.
- $x = 2$ implies $x^2 = 4$.
- $x^2 = 4$ if $x = 2$.
- $x = 2 \Rightarrow x^2 = 4$.

Which form to use is a matter of style and, most importantly, clarity. Notice how the arrow symbol in the last example so eloquently conveys the condition $x = 2$ leading to the consequence $x^2 = 4$. On the other hand, for the more complicated proposition that if $x = x'$ and $y' = y'$, then $x + x' = y + y'$, the form “ $x = x'$ and $y' = y' \Rightarrow x + x' = y + y'$ ” might be ambiguous, appearing to possibly mean that $x = x'$ (unconditionally) and, in addition, if $y = y'$, then $x + x' = y + y'$, an entirely different meaning, and one that is certainly not universally true for all real numbers x , x' , y , and y' !

The *converse* of an implication is the proposition obtained by exchanging the condition and the consequence. The converse of a true proposition need not be true. For example, the converse of the universally true proposition that if $x > 2$, then $x^2 > 4$, is the proposition that if $x^2 > 4$, then $x > 2$, which fails to be universally true. (The latter proposition fails because a number that satisfies the condition $x^2 > 4$ could be less than -2 ; the condition is not sufficient to ensure the stated consequence in all instances.) We often summarize the conjunction of an implication and its converse with the phrase “if and only if” or an evocative double arrow, as in the following (universally true) logically equivalent statements:

- $x + 2 = 5$ if and only if $x = 3$.
- $x = 3$ if and only if $x + 2 = 5$.
- $x + 2 = 5 \Leftrightarrow x = 3$.
- $x = 3 \Leftrightarrow x + 2 = 5$.

Suppose $p(x)$ and $q(x)$ are open propositions about a real number x for which the assertion $p(x) \Leftrightarrow q(x)$ holds universally (that is, for all $x \in \mathbb{R}$). For any instance $x \in \mathbb{R}$ such that $p(x)$ is true, $q(x)$ must also be true (since $p(x) \Rightarrow q(x)$). Similarly, for any instance $x \in \mathbb{R}$ such that $q(x)$ is true, $p(x)$ must also be true (since $q(x) \Rightarrow p(x)$). Thus, $p(x)$ and $q(x)$ are true for exactly the same instances and false for exactly the same instances: they are equivalent.

Turning now to the remaining assertions, the associative property of addition is also a universal property: for all real numbers x, y , and z , $(x+y)+z = x+(y+z)$. (Students occasionally confuse the associative property with the commutative property, which is also assumed to be true for addition in any number system. The commutative property involves only two terms, stating that for all real numbers x and y , $x+y = y+x$. The associative property involves three terms, and the order of the terms is the same on both sides of the equation; it is just a matter of which pair of adjacent terms is summed first, after which the third term is added in whatever order it appears. The associative property simply states that the result of successive additions does not depend on how the terms are successively grouped into pairs.)

The defining property of 0 has both universal and existential aspects: First we assert that there *exists* a real number y such that, for *every* real number x , $x+y = y$. Notice that we are asserting the *existence* of a number with a particular *universal* property. For convenience, we might refer to this property as the identity property under addition: the result of adding this number to any other given number is the identical given number you started with. It is easy to deduce that there is only one number with this property, for suppose x and x' both have it. Then $x = x + x' = x' + x = x'$. The first equation holds by the identity property assumed for x' , the second follows from the commutative property of addition, and the third from the identity property assumed for x . Since there is only one such number, we can give it a name: 0. In summary, we define 0 to be the unique number with the property that, for any real number x , $x+0 = 0$. The number 0 is called the *additive identity*. (By the commutative property of addition, we also have that for any real number x , $0+x = 0$. For simplicity, we include this fact in the defining property of 0. That way, we don't have to remember in which order we stated it and apply the commutative property every time 0 appears in the other position!)

Similarly, the existence of additive inverses includes both universal and existential aspects, but applied in the opposite order: for *any* real number x , there *exists* a real number y such that $x+y = 0$. Here we are asserting the *universal existence* of numbers having properties specific to the instance to which they apply. We will prove later on that each real number x has a unique additive inverse, which we denote by $-x$. Our assumption is only that it has an additive inverse, not that it has only one. (You can probably figure out the argument now. Try it!)

Now let us return to what seemed to be a trivial equation, knowing much more about what its statement and solution embody! Although you may remember solving equations as a sequence of manipulative “bookkeeping” steps, these steps represent a process of logical deduction in which properties such as these are applied. If a number with a particular property, such as 0, -2 , or $\sqrt{2}$ exists, then we can introduce this number as appropriate and take advantage of its defining property. If a conditional property applies universally, we can use it in any instance for which the condition is satisfied to obtain the conclusion. Unfortunately, in many elementary treatments of algebra the logical relationships are omitted from the written manipulations. Let us put them in. In the context of the real number system:

- $x + 2 = 5 \Leftrightarrow (x + 2) + (-2) = 5 + (-2)$, since addition is a well defined operation (operative quality of addition). The converse implication (\Leftarrow) holds because we could undo what we did by adding 2 to both sides of our new equation, retrieving the original one. Our introduction of the number (-2) is justified by the assumption that a unique additive inverse exists for any real number, along with the definition of -2 as denoting the additive inverse of 2.
- $(x + 2) + (-2) = x + (2 + (-2))$, by application of the associative property of addition. As with the operative quality of addition, we do not need to know what x is to apply this property, since it is universal.
- $2 + (-2) = 0$ by the defining property of -2 , and $x + 0 = x$, by the defining property of 0; therefore, $x + (2 + (-2)) = x$.
- On the other hand, $5 + (-2) = 3$. This calculation follows from the definitions of 5, 4, 2, -2 , and 0, along with the associative property of addition: $5 + (-2) = (4 + 1) + (-2) = ((3 + 1) + 1) + (-2) = (3 + (1 + 1)) + (-2) = (3 + 2) + (-2) = 3 + (2 + (-2)) = 3 + 0 = 3!$
- Thus, $x + 2 = 5 \Leftrightarrow x = 3$.

Since writing all of this down is rather cumbersome and buries the actual process among the justifications, we generally leave commonplace justifications to the reader when writing the solution to an equation. (Because the associative property is used so routinely, we usually don’t even bother to group terms when the only operation is addition.) But we should at least retain the logical relationships between the equations, as in the following streamlined exposition, rather than just writing them in sequence:

$$x + 2 = 5 \Leftrightarrow x + 2 + (-2) = 5 + (-2) \Leftrightarrow x = 3.$$

I hope you will do this from now on. The value of this practice will be more apparent in the solution of a more complicated equation, which we address after some exercises.

3.3 Exercises

1. Write each of the following sets in formal notation.
 - (a) The set of real numbers whose squares are greater than 1 and less than 3.
 - (b) The set of integers whose squares are greater than 1 and less than 3.
2.
 - (a) Write the converse of the proposition that, if $x < 5$, then $x - 1 < 5$.
 - (b) Is the proposition you just constructed true for all real numbers x ?
 - (c) Is it existentially true? That is, do there exist real numbers x for which it is true?
 - (d) If it is existentially true, is this fact meaningful in the sense of providing new information and predictive power?
3.
 - (a) Write the converse of the proposition that if $x = x'$ and $y = y'$, then $x + x' = y + y'$.
 - (b) Is the proposition you just constructed true for all real numbers x, x', y , and y' ?
 - (c) Is it existentially true?
 - (d) If it is existentially true, is this fact meaningful in the sense of providing new information and predictive power?
4. Comment in general on whether existentially true implications have argumentative power.
5. For which of the following sets is the proposition that $x^2 > 4 \Rightarrow x > 2$ universally true?
 - \mathbb{N}
 - $\{x \in \mathbb{R} : x > 0\}$
 - $\{x \in \mathbb{R} : x > -2\}$
 - $\{x \in \mathbb{R} : x > -1\}$

6. If we think of conjunction, disjunction, and implication as binary logical operations, we can ask whether each is associative or commutative.
 - (a) Which of these operations are associative?
 - (b) If any are not associative, give an example of an open proposition whose truth values change when its components are regrouped.
 - (c) Which of these operations are commutative?
 - (d) If any are not commutative, give an example of an open proposition whose truth values change when its components are exchanged.
7.
 - (a) What property would a multiplicative identity have?
 - (b) Does a number having this property exist?
 - (c) Must a number with this property be unique? If so, prove it!
 - (d) If a unique such number exists, what is it?
8. For each of the following propositions, state if it is true or false, and briefly sketch an argument to justify your answer. Remember that, as we saw in the examples in the text, the syntactical convention is that anything asserted to exist may depend on those things, and only on those things, that have been introduced earlier in the proposition.
 - (a) For any real number y , there exists a real number x such that $x < y$.
 - (b) There exists a real number x such that, for any real number y , $x < y$.

3.4 A More Complicated Equation

Let us now find a list of all real solutions to the equation $x = \sqrt{x+2}$, providing logical justification that our solution is correct. (Recall that the symbol “ $\sqrt{\quad}$ ” denotes the *non-negative* square root, thereby designating a unique number.) In doing so, we will refer to some familiar properties assumed for numerical operations that we have not listed before. We will also need the following theorem, which you will prove rigorously in Chapter 5: For all real numbers x and y , $xy = 0$ if and only if $x = 0$ or $y = 0$. For now, here is a sketch of the proof. First, to prove the implication from right to left, we show that, for any real number x , $0x = 0$. (Note that the defining property of 0 does not say this! Since 0 has already be defined to satisfy a different property, we cannot simply assume it also has this one: we have to prove it does!) $0x = (0 + 0)x$ by the defining

property of 0, and $(0 + 0)x = 0x + 0x$ by the distributive property; hence, $0x = 0x + 0x$. Subtracting $0x$ from both sides, we obtain the result. Conversely, suppose that $xy = 0$, but $x \neq 0$; in this case, we must show that $y = 0$. Since $x \neq 0$, we can divide both sides by x to conclude that indeed $y = 0 \cdot \frac{1}{x} = 0$.

Now, on to our solution to the equation:

- For any real number x , $x = \sqrt{x+2} \Rightarrow x^2 = x+2$ by the operative quality of multiplication: on the left, we have multiplied a number by itself; on the right, we have multiplied the same number (albeit expressed differently) by itself; same inputs, same output.

Observe that the converse implication is *not* universally true. Why not? Since we do not know the value of x , we cannot apply any property that is not universal to all real numbers.

- For any real number x , $x^2 = x+2 \Leftrightarrow x^2 - x - 2 = 0$. Justification is left to the reader.
- By the distributive property, $x^2 - x - 2 = 0 \Leftrightarrow (x-2)(x+1) = 0$.
- By the theorem cited above, $(x-2)(x+1) = 0$ if and only if $x-2 = 0$ or $x+1 = 0$.
- $x-2 = 0 \Leftrightarrow x = 2$; similarly, $x+1 = 0 \Leftrightarrow x = -1$. Justification is left to the reader.
- Thus, putting together our previous steps, we obtain that if $x = \sqrt{x+2}$, then $x = 2$ or $x = -1$. (We do not obtain the converse implication because the first step is not “reversible.”) This implication shows that 2 and -1 are the only possible solutions to the equation $x = \sqrt{x+2}$, but it does not guarantee that each of these numbers is a solution. (The converse implication, which fails, would guarantee they were if it were true.) Therefore we check each solution, finding that indeed 2 solves the equation, but that -1 does not (for reasons that should be obvious).
- We conclude that $\{x \in \mathbb{R} : x = \sqrt{x+2}\} = \{2\}$.

3.5 Exercises

1. Explain why it is not universally true for all real numbers that if $x^2 = x+2$, then $x = \sqrt{x+2}$.

2. Describe a useful set of real numbers for which it *is* universally true that if $x^2 = x + 2$, then $x = \sqrt{x + 2}$. (Don't use what you know about the solution to this particular equation, since then the set you come up with will not be useful for solving it!)
3. Based on your answer to the previous question, is it necessary from a logical point of view (as opposed to the possibility of having made an error in calculation) to check that 2 is a solution to the equation $x = \sqrt{x + 2}$? (In other words, can we tell immediately by logical deduction that, as long as we have not made an error in calculation, 2 must be a solution.)
4. Is the proposition that $x \neq \sqrt{x + 2}$ universally true for the set $\{x \in \mathbb{R} : x < 0\}$?
5. Based on your answer to the previous question, can we tell immediately without checking by calculation that -1 is not a solution to the equation $x = \sqrt{x + 2}$?
6. Justify the proposition that, for any real number x , $x^2 = x + 2 \Leftrightarrow x^2 - x - 2 = 0$.

Chapter 4

Sets, Logic, & Proof

In which we systematically summarize the the relationship between sets and logic, introduce some important logical relationships, and consider the structure of a proof. In doing so, we feel impelled to venture into foundational and philosophical areas, but not too much!

I hope the previous chapters have convinced you that you have defined sets and made logical deductions more often than you may have thought, even if you did not precisely formulate your definitions and arguments and were not fully conscious of the extent to which sets and logic underlie elementary geometric and algebraic procedures! Now we begin a rigorous treatment of these topics.

4.1 Open Propositions Define Sets

Just as an equation defines its solution set within a given number system, any open proposition whose truth value depends on a single variable defines the set of all elements of a given set that make the propositions true, as in the following examples:

- $\{x \in \mathbb{R} : x = \sqrt{x+2}\} = \{2\}$
- $\{x \in \mathbb{N} : 2 < x < 5\} = \{3, 4\}$
- $\{x \in \mathbb{Q} : 2 < x < 5\} = \{3, 4, \frac{5}{2}, \frac{7}{2}, \frac{9}{2}, \frac{7}{3}, \frac{8}{3}, \frac{10}{3}, \dots\}$
- The elements of $\{x \in \mathbb{R} : 2 < x < 5\}$ cannot be listed, even with an infinite list.
- $\{x \in \mathbb{N} : \text{For every } y \in \mathbb{N}, x \leq y\} = \{0\}$

- $\{x \in \mathbb{R} : \text{For every } y \in \mathbb{N}, x \leq y\} = \{x \in \mathbb{R} : x \leq 0\}$
- $\{x \in \mathbb{N} : \text{For every } y \in \mathbb{R}, x \leq y\} = \{x \in \mathbb{R} : \text{For every } y \in \mathbb{R}, x \leq y\} = \emptyset$

The proposition that defines a set is said to *separate* the elements that satisfy it from the remaining elements of the set initially given. The set defined by the proposition is said to *comprehend* it - in more common usage, it is the comprehensive collection of elements that satisfy the proposition - subject to the *restriction* that its elements come from the given set. To comprehend means to understand, and in some sense, at the most basic level, to understand a property means to know all the things that have it, the set determined by what it means. (Certainly, there are higher levels of understanding as well: seeing parallels, analogies, and metaphors, for example.) Of course, as illustrated by the examples above, the initial set, by restricting the “universe” from which elements are selected, makes a considerable difference to the set that is defined.

There are a few occasions when, discussing very general considerations, we cannot specify an initial set to which all elements of a set we define must belong. For example, suppose we are given two sets A and B whose elements are not specified. For any two sets A and B , it certainly makes sense to define the set of all elements that belong either to A or to B (possibly both, as we never use “or” in the exclusive sense without specifically saying so). The resulting set, called the *union* of A and B , depends only on the sets A and B :

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

Since we wish to make this definition in general, for whatever sets A and B might arise, we have no way of choosing an initial set to which the potential elements of $A \cup B$ should be restricted.

There are two other important operations defined on sets:

Definition. Given sets A and B , $A \cap B = \{x : x \in A \text{ and } x \in B\}$

The set $A \cap B$ is called the *intersection* of A and B .

Definition. Given sets A and B , $A \setminus B = \{x : x \in A \text{ and } x \notin B\}$

The set $A \setminus B$ is called the *complement* of B in A . The set $(A \setminus B) \cup (B \setminus A)$ is called the *symmetric difference* of A and B , commonly denoted by $A \triangle B$ or $A \ominus B$.

The operations of union and intersection can be defined more generally, which is necessary in order to consider the union or intersection of infinitely many sets. If we could only perform these operations on two sets at a time, we could only create the

union of finitely many sets, no matter how many times we iterated this process. (The term “finite” can be precisely defined, but we will be satisfied with its intuitive meaning for purposes of this discussion.) The following definitions generalize these operations to any collection of sets.

Definition. Given a set \mathcal{A} (whose elements are assumed also to be sets),

$$\bigcup_{A \in \mathcal{A}} A = \{x : x \in A \text{ for some } A \in \mathcal{A}\}.$$

Definition. Given a set \mathcal{A} (whose elements are assumed also to be sets),

$$\bigcap_{A \in \mathcal{A}} A = \{x : x \in A \text{ for every } A \in \mathcal{A}\}.$$

Until about a century ago, mathematicians did not worry about making general definitions of sets, such as the ones we just made, whenever it seemed convenient to do so. But then, in the early twentieth century, the philosopher and mathematician Bertrand Russell discovered a danger. Recall that we informally described a set, not merely as a collection of things, but as an *identifiable* collection of things. What this description means to suggest is that, for any mathematical object, the question of whether or not this object is in a given set should have an answer. We might not know the answer. We might not even know a procedure for determining an answer. But surely, for any reasonable conceptualization of sets and mathematical objects, the question should have an answer, and this answer should be exactly one of two choices: “yes” or “no”! Formally, given any defined set A and any mathematical object x , $x \in A$ should be a closed proposition: it is either true or false, and not both. What Russell discovered was a purported definition for a set A , along with an indisputably mathematical object x for which the statement $x \in A$ does not have a unique truth value! (We say a “purported” definition because, if it cannot be determined for all objects x whether or not x belongs to the set, then surely we cannot say the set has actually been defined at all.)

Before exhibiting Russell’s troubling example, let us make Russell’s conceptualization of sets and mathematical objects more concrete: mathematical objects and sets are, simply, the same thing. It should come as no surprise that sets are mathematical objects. Indeed, if mathematics is the study of patterns, and if patterns are recognized by distinguishing and classifying their parts, and if classifying things means putting them into categories, and if defining categories is the same as defining the sets of things in these categories, then surely sets are mathematical objects! On the other hand, it may seem a bit peculiar at first that the *only* mathematical objects are sets. After all, you are used to working with such things as numbers and geometric points. Are these things sets? According to this point of view, yes, they are. That is, we can define sets that represent

them, we can define the operations that apply to them as set operations, and we can define the relationships between them as relationships between sets. In particular, there is no reason to shy away from letting sets be elements of other sets, and indeed those are the only kind of elements there are. We develop this perspective further in Chapter ??.

Note that a set of sets is not the same as the union of these sets, as illustrated by the following examples:

- $\{\{1, 2\}, \{2, 3\}\}$ has two elements, but $\{1, 2\} \cup \{2, 3\}$ has three elements.
- $\{\{1, 2, 3\}, \{2, 3, 4, 5, 6\}\}$ has two elements, but $\{1, 2, 3\} \cup \{2, 3, 4, 5, 6\}$ has six elements.
- $\{\mathbb{N}, \mathbb{Q}\}$ has two elements, but $\mathbb{N} \cup \mathbb{Q} = \mathbb{Q}$ has infinitely many elements, which can be listed systematically.
- $\{\mathbb{N}, \mathbb{Q}, \mathbb{R}\}$ has three elements, but $(\mathbb{N} \cup \mathbb{Q}) \cup \mathbb{R} = \mathbb{Q} \cup \mathbb{R} = \mathbb{R}$ has infinitely many elements, which cannot be listed.

Here is Russell's example. Suppose we attempt to define a set A in the following manner:

$$A = \{x : x \notin x\}$$

At first this proposed definition might seem reasonable: if X is a specific set, then give sufficient information about the elements of X , such as that given by its definition, we should be able to check whether or not X itself is an element of X . But consider the set A specified by the proposed definition. If this is a good definition, then A is a set (and hence an object for mathematical consideration). If A is a set, then we must be able to tell if A itself is a member of A : the statement that $A \in A$ must be a closed proposition that is uniquely either true or false. Suppose this statement is true. Then, by definition of A , it is false! On the other hand, suppose the statement that $A \in A$ is false. In this case, $A \notin A$ is true; hence by definition of A , the statement that $A \in A$ is true! Clearly, the statement that $A \in A$ does not have a unique truth value. Thus, it is not a proposition. So A cannot be a set!

If a proposed definition does not actually ensure that what it purports to define, whether a quality or an object, can be identified (given sufficient information), then it must be rejected. We say that what it purports to define is not *well* defined. The set A of Russell's example is not well defined.

In response to this problem, mathematicians have devised specific assumptions, or axioms, about what sets (and hence what "mathematical objects") exist. Any assertion

we wish to make about sets, including the definition of any specific set (since to define it is to assert that it exists), must ultimately be justified by the axioms of set theory (of which there are several variations). We may assume, as an axiom of some other theory, that a set with some specific properties exists; we will in fact do so in the following chapter for the set of real numbers and the operations of addition and multiplication. However, any such assumption is provisional. That such a set actually does exist must be justified from the axioms of set theory.

There certainly can be nothing wrong with the open statement “ $x \notin x$ ” itself, which is constructed using only the membership relation and the logical operation of negation. The problem is that it has been applied in an unrestricted manner, beyond the instances for which its meaning can be comprehended. The axiom that prevents Russell’s paradox, called the Axiom of Restricted Comprehension, restricts the application of a property to the elements of an existing set. (The Axiom of Restricted Comprehension is also called the Axiom of Separation, among other things, since it separates the elements that satisfy it from those that don’t.)

Axiom (Axiom of Restricted Comprehension). *For all sets A , if $p(x)$ is an open proposition whose truth value depends only on the value of the variable x , then $\{x \in A : p(x)\}$ is a set.*

In reality, the Axiom of Restricted Comprehension is a whole family of axioms, called an axiom schema, one for each open proposition $p(x)$: since propositions are linguistic constructions about mathematical objects, not mathematical objects themselves, we cannot refer to “all open propositions” in a proposition! (In other words, we have provided a scheme for creating all of the axioms in this family: just pick any open proposition whose truth value depends on a single variable and plug it in for $p(x)$. There is an explicit process for constructing any open proposition in the formal language of set theory, using only variables, the relations $=$ and \in , and logical operations, but we have no need for that level of formality at this point.) Note that $\{x \in A : p(x)\}$ is synonymous with $\{x : x \in A \text{ and } p(x)\}$. Thus, all of the axioms ultimately involve a variable whose values are unrestricted. (The Axiom of Restricted Comprehension says that if this variable is restricted to be a member of a known set in conjunction with satisfying some other proposition, then the resulting set is well defined.)

The Axiom of Restricted Comprehension creates some new problems for us. For one thing, in order to define any sets at all, we must first have a set! For another, even if we did have some sets, our definition of the union operation is no longer legitimate, since we don’t have a set to which the proposition that $x \in A$ or $x \in B$, or the more general proposition that for some $A \in \mathcal{A}$, $x \in A$, should be restricted. (The union $A \cup B$, or more generally $\bigcup_{A \in \mathcal{A}} A$, would be such a set, but we cannot define it until we already

have defined it, a circular conundrum!) There is only one way out: we need some axioms ensuring the existence of some sets.

Our first axiom ensures the existence of the most basic set of all: the empty set. This axiom is generally not included separately in the most generally used standard axioms of set theory, as it follows in a rather subtle way from a much stronger, but harder to understand, axiom that we will leave for Chapter ??; however, for simplicity and clarity, we include it.

Axiom (Existence of the Empty Set). *There is a set with no elements. That is, there exists a set \emptyset for which the following (unrestricted) proposition is true: for every x , $x \notin \emptyset$.*

The next axiom allows us to create sets with up to two elements - still pretty basic.

Axiom (Axiom of Pairing). *For all sets A and B , $\{A, B\}$ is a set.*

By $\{A, B\}$ we mean, of course, the set for which the proposition that $x \in \{A, B\}$ is true for the instances $x = A$ and $x = B$, false in all other instances. Note that the sets A and B need not be distinct. So, for example, we can now justify that $\{\emptyset\}$ is a set.

Axiom (Axiom of Union). *For any set \mathcal{A} (whose elements are assumed to be sets),*

$\bigcup_{A \in \mathcal{A}} A = \{x : x \in A \text{ for some } A \in \mathcal{A}\}$ *is a set.*

Note that together with the Axiom of Pairing, the Axiom of Union provides in particular for the union of two sets, since $A \cup B = \bigcup_{X \in \{A, B\}} X$. Starting with the empty set and using the Axioms of Pairing and Union, we can create sets with any finite number of elements, but no infinite set. And, of course, any set separated from a finite set by an open proposition will be finite. Thus, these four axioms are not sufficient to justify the existence of even the natural numbers. Surely we cannot have mathematics without even being able to count indefinitely! Of course not! But we will leave the the axiom that asserts the existence of the natural numbers - the much stronger but harder to understand axiom mentioned above - for Chapter ??.

You might reasonably ask if the axioms of set theory succeed in solving the problem illuminated by Russell, not only avoiding Russell's specific paradox (which you will show in the exercises that they do), but preventing us in general from ever defining a set that does not make logical sense. The answer is that we don't know, and in fact there is some controversy among mathematicians about whether or not it is safe to assume one of them in particular (which we have not yet introduced - the one's introduced so far are uncontroversial). If someone were to find some set defined in accordance with the axioms that leads to a logical contradiction, as Russell's did, and hence is not well defined, then

we would know the axioms were not adequate. We would have to discard some of the axioms and perhaps formulate some new ones and hope that they were better, and we would have to discard any mathematics that was based on the faulty definition. So far, no one has done this, but we cannot guarantee that no one ever will. Mathematics is never finished!

Fortunately, people do seem to have good instincts about the nature of mathematics. Mathematicians certainly make mistakes, but so far in the history of mathematics no serious mathematics has ever been discarded because it exposes a fundamental flaw in the logical foundation of the subject. Although Russell's reasoning certainly counts as insightful and serious mathematics, discarding his ill-defined example does not count as a loss to mathematics, because Russell's only reason in formulating it was to expose the vulnerabilities in the logical foundation being used at that time, much as a computer hacker might expose the vulnerabilities in the security of a network. He showed that the set theory of the time was naive and in need of further development. (Indeed we still call treatments of set operations without reference to foundational issues of set definition "naive" set theory.) Fortunately, Russell exposed these vulnerabilities before anyone made a serious breach of logic (at least as far as we know).

4.2 Exercises

1. Let B be an existing set. Show that no contradiction arises from the restricted definition $A = \{x \in B : x \notin x\}$, and determine whether or not $A \in A$. Thus, the Axiom of Restricted Comprehension avoids Russell's paradox.
2. For sets A and B as in the previous exercise, determine whether or not $B \in A$.

4.3 Logical Equivalence

Different combinations of operations may produce the same result for all inputs. As a simple example from algebra, for any real number x , $2x = x + x$. Although multiplying a number by 2 is operationally different from adding a number to itself, the properties of these operations ensure that they always give the same result: by definition, $2 = 1 + 1$; the distributive property implies that $(1 + 1)x = 1x + 1x$; and the multiplicative identity property of 1 implies that $1x + 1x = x + x$. Thus, we might say that these different procedures are *operatively equivalent* and that the expressions " $2x$ " and " $x + x$ " (as opposed to the numbers they represent, which are equal) are *expressively equivalent*.

Just as an algebraic expression represents a number for each instance of its variables, a proposition represents a truth value in each instance. As we discussed in Chapter 3, propositions with the same truth value in each instance are considered *equivalent*. Since equivalent propositions define the same set when one is substituted for the other, they have the same categorical value. (They may well have different computational value. A large part of mathematics consists of finding equivalent criteria for that are easier to compute.) If their equivalence derives solely from the properties of the logical operations used, rather than the content of the component propositions involved, the propositions are said to be *logically equivalent*. Logically equivalent forms yield equivalent propositions for all sets of propositions to which they are applied. Just as the expressive equivalence of “ $2x$ ” and “ $x + x$ ” has nothing to do with a particular value of x , but rather derives solely from the way these expressions are constructed, logical equivalence has nothing to do with the meaning of the particular propositions that are combined, but results entirely from the way they are combined, from the form rather than the content. The propositions “ $x \not\leq y$ ” and “ $x \geq y$ ” are equivalent, but not *logically* equivalent, since their equivalence derives from the properties of the particular relation of order in the real number system. (Don’t let the terminology confuse you. It is not that “ $x \not\leq y$ ” and “ $x \geq y$ ” are not equivalent for all logical purposes - that is, for purposes of argument. They are, at least if x and y are real numbers and $<$ denotes the usual order in the real number system. It is just that they are not equivalent for reasons that are purely intrinsic to the nature of logical expression. Additional assumptions are required.) Let us now consider some very important examples of logically equivalent forms that propositions may take.

4.3.1 Contrapositives

Consider the propositions “ $x > 2 \Rightarrow x^2 > 4$ ” and “ $x^2 \not> 4 \Rightarrow x \not> 2$.” These two propositions represent the general forms “ $p \Rightarrow q$ ” and “ $\neg q \Rightarrow \neg p$,” respectively, with the particular propositions “ $x > 2$ ” and “ $x^2 > 4$ ” substituted for p and q , respectively. Each of these conditional propositions is called the *contrapositive* of the other. (Why the term is not “contranegative” is anyone’s guess, but that’s the way it is!) No matter what propositions are substituted for p and q , contrapositives have the same meaning: stating that condition p ensures consequence q is the same as stating that the failure of said consequence q ensures the failure of said condition p (since if p had not failed, q could not have failed, whereas the satisfaction of q tells us nothing about p).

Different equivalent expressions are useful in different situations. For example, of the equivalent expressions “ $x^2 - 4$ ” and “ $(x + 2)(x - 2)$,” the former is more useful when computing the derivative of the function f defined by $f(x) = x^2 - 4$, whereas the latter is more useful when finding the roots of $f(x) = 0$. Similarly, one of two contrapositively

related propositions might be more useful in a particular proof, because these equivalent propositions provide different hypotheses from which to work.

4.3.2 Rephrasing a conditional as a disjunction

At first glance, it might not seem that one could rephrase a conditional proposition as a disjunction, or vice-versa, because we tend to use these constructions in different senses. We generally use conditional propositions in a *predictive* sense. Suppose I say, “If you get a C or better on the third exam, you will pass the course.” Then, presuming that I am not lying, when you get your third exam back having received, say, a B, you can predict with certainty that you won’t fail the course. If this course was the last obstacle to graduation, you can safely tell your parents to make their hotel reservations! On the other hand, we typically use disjunctive propositions *descriptively*. When I say, “Either you got a D or F on the third exam or you passed the course,” I am simply describing all the possibilities that might have occurred.

But notice that if this disjunction is stated in the future tense - “Either you will get a D or F on the third exam or you will pass the course” - it has exactly the same predictive value as the previous conditional proposition. If those are the only two possibilities, and at least one must be true, then clearly if you get a C or better (and hence do not get a D or F), you will pass. It is also possible that both are true. Maybe I am being conservative in my prediction or just pushing you to get at least a C on the exam; perhaps you could get a D and still pass. (In everyday usage, an unstated converse implication is sometimes implicitly also intended, but in mathematics we do not recognize any propositions that are not explicitly stated!) Similarly, if the previous conditional statement is stated in the past tense - “If you got a C or better on the third exam, then you passed the course” - it takes on a descriptive sense, narrowing the possibilities to having either got less than a C on the third exam or passed the course, one of which must have occurred. It is also worth noting at this point that both propositions *exclude* exactly the same single possibility: getting a C or better on the third exam and failing the course. (In everyday speech, the use of a disjunction in the predictive sense seems to be used most often for emphasis when the consequence is negative, as in “You will be home by midnight or you will be grounded for two weeks”!)

From a logical point of view, sense is irrelevant; only truth value is relevant. Thus the forms “ $p \Rightarrow q$ ” and “ $\neg p$ or q ” are logically equivalent, as I hope the preceding discussion has convinced you. The assertion of each of these forms excludes the possibility that p holds and q does not, while allowing the other three (that both hold, that both fail, or that p fails and q does not). I also hope that by informally considering the meaning of examples you are getting a feel for how logical language is used. The cartoons in Figure

4.1 may further help make the point in a darkly humorous manner!

As with contrapositives, the disjunctive and conditional forms have different uses. As we will see shortly, a disjunction provides a natural division into cases (see Section 4.6.4), whereas the conditional form provides a hypothesis on which to rest subsequent arguments.

4.3.3 Negation: Don't just say no!

Often we must assert that something is not the case; usually it is more useful to do so in an affirmative manner, stating exactly what *is* then the case. A note on syntactical conventions: the phrase “it is not the case” is assumed to apply to all that follows before the next punctuation mark. So, “It is not the case that p and q ” has the same meaning as “ $\neg(p \text{ and } q)$ ”, whereas “It is not the case that p , and q ” has the same meaning as “ $(\neg p)$ and q .” In spite of these conventions, negative phrasing can be confusing, another reason to rephrase the proposition with all the negatives (if any) in the atomic component propositions. Our understanding of truth value shows us how to do this.

Considering a real number x , suppose it is not the case that $x \not> 2$. Then clearly it is the case that $x > 2$. For any proposition, the negation of its negation is logically equivalent to the original proposition. Symbolically, “ $\neg(\neg p)$ ” is logically equivalent to “ p .”

Again considering a real number x , suppose it is not the case that $x > 2$ or $x \in \mathbb{Z}$; then it is the case that $x \not> 2$ and $x \notin \mathbb{Z}$. For any real number x , there are exactly four mutually exclusive possibilities: $x > 2$ and $x \in \mathbb{Z}$, $x > 2$ and $x \notin \mathbb{Z}$, $x \not> 2$ and $x \in \mathbb{Z}$, or $x \not> 2$ and $x \notin \mathbb{Z}$. The assertion that $x > 2$ or $x \in \mathbb{Z}$ includes the first three cases and excludes the fourth; the negation of this assertion, then, includes exactly the fourth: $x \not> 2$ and $x \notin \mathbb{Z}$. In general, given any two propositions p and q , there are, *a priori*, exactly four mutually exclusive possibilities: p and q , p and $(\neg q)$, $(\neg p)$ and q , or $(\neg p)$ and $(\neg q)$. (I say “*a priori*” because a mathematical relationship between particular propositions p and q might preclude some of them. For example, the case that $x > 2$ and $x \not> 1$ does not occur.) The assertion “ p or q ” includes the first three and excludes the fourth. So if it is not the case that p or q , then it is the case that $(\neg p)$ and $(\neg q)$. Symbolically, “ $\neg(p \text{ or } q)$ ” is logically equivalent to “ $(\neg p)$ and $(\neg q)$.”

Since “ $\neg(\neg p)$ ” is logically equivalent to “ p ,” it is now easy to see, working backwards, that the negation of “ p and q ” is “ $(\neg p)$ or $(\neg q)$.” Since “ $p \Rightarrow q$ ” is logically equivalent to “ $(\neg p)$ or q ,” its negation is logically equivalent to “ p and $(\neg q)$.” If we are thinking of the assertion “ $p \Rightarrow q$ ” in the predictive sense, then its negation “ p and $(\neg q)$ ” asserts the case in which this predictive device does not work.

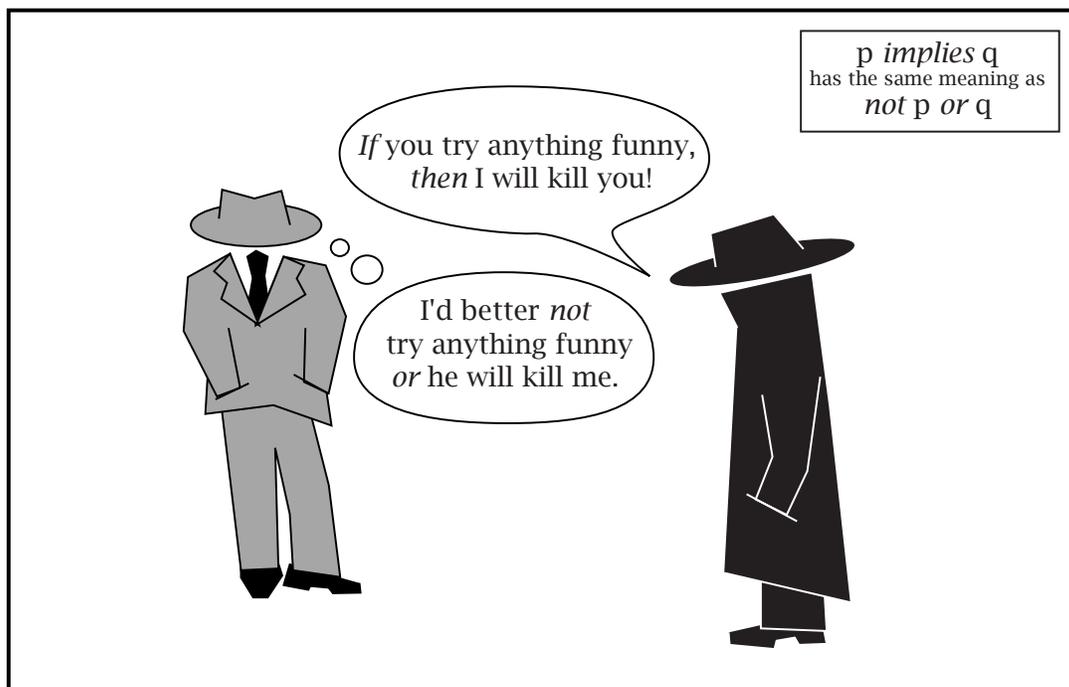
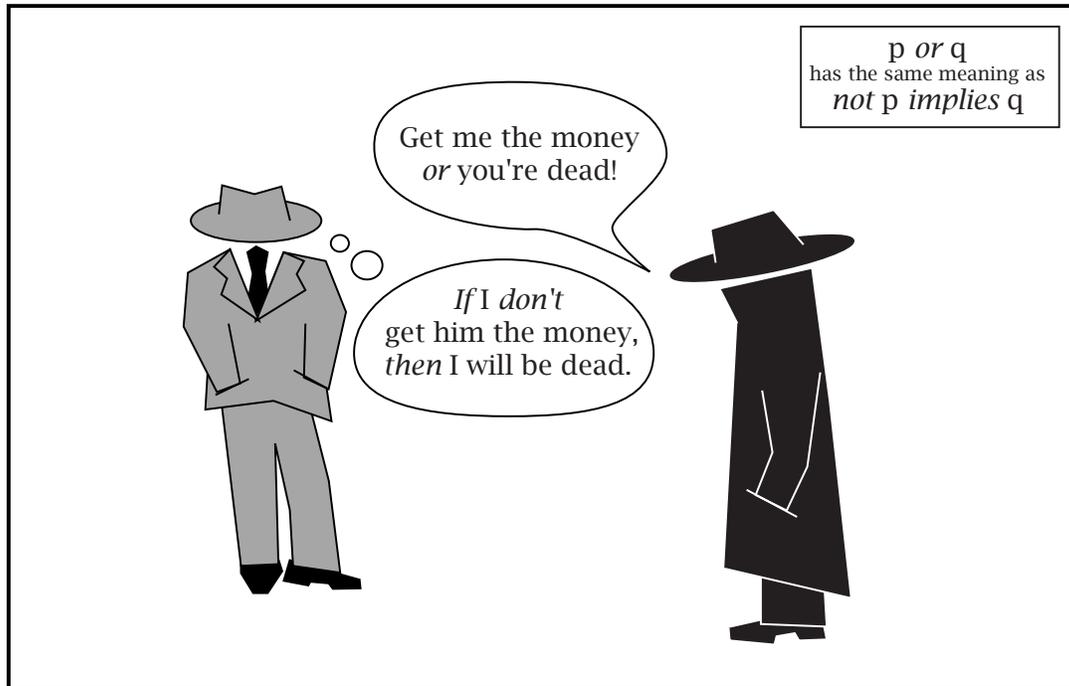


Figure 4.1: Logic on the dark side!

Example. If it is not the case that $2 < x < 5$ (shorthand for $2 < x$ and $x < 5$), then it is the case that $2 \not< x$ or $x \not< 5$.

Example. If it is not the case that $x^2 = x + 2 \Rightarrow x = \sqrt{x + 2}$, then it is the case $x^2 = x + 2$ and $x \neq \sqrt{x + 2}$.

Example. The proposition that $x \notin A \cup B$ means it is not the case that $x \in A \cup B$, which by definition means it is not the case that $x \in A$ or $x \in B$. Rephrasing this proposition in a more useful, affirmative manner, with all the negations moved into the atomic components, we have that $x \notin A$ and $x \notin B$.

Example. The proposition that $x \notin A \cap B$ means it is not the case that $x \in A \cap B$, which by definition means it is not the case that $x \in A$ and $x \in B$. Rephrasing this proposition in a more useful, affirmative manner, with all the negations moved into the atomic components, we have that $x \notin A$ or $x \notin B$.

For complicated propositions, we just work our way in:

Example. Suppose it is not the case that if $2 < x < 5$, then $x^2 < 17$ or $x^2 > 8$. Then it is the case that $2 < x < 5$ and $x^2 \not< 17$ and $x^2 \not> 8$.

Remark. If x is confined to be an integer, then it is the case that if $2 < x < 5$, then $x^2 < 17$ or $x^2 > 8$.

Let us pause to observe a fundamental difference in the nature of conjunction, on the one hand, and that of disjunction or implication, on the other. A conjunction of closed propositions whose truth value we know provides a useful summary by asserting two facts together in a neat package, as in $2 < 5$ and $5 < 7$. (Informally, 5 lies between 2 and 7.) But a disjunction of two propositions whose truth value we know is not really useful, because it asserts less than we really know. Although $2 < 5$ or $5 < 7$ is certainly true, it would be better to assert the stronger proposition stated above. And although $2 < 5$ or $5 < 2$ is also certainly true, it would be better simply to assert the stronger proposition that $2 < 5$; throwing $5 < 2$, which we know to be false, into the mix is a distraction that results in a loss of information. Rephrasing either of these disjunctions as an implication also seems silly: in the statement $2 \not< 5 \Rightarrow 5 < 7$, the condition is meaningless, since the conclusion that $5 < 7$ is unconditionally true, and in the statement $2 \not< 5 \Rightarrow 5 < 2$, the false conclusion is meaningless, since the condition is patently false. Disjunctions and implications only come into their own when they apply to propositions whose truth value we don't know, because in this context they define relationships between these unknowns. In particular, disjunctions and implications of open propositions are valuable

when they apply universally, either to all mathematical objects or, more commonly, to all the elements of some set. (They are most useful when the collection of objects to which they apply is not finite, since for a finite collection we can check each instance, at least in theory.) Thus, we come naturally to the subject of *quantification*.

4.4 Quantified Propositions Define Relationships

4.4.1 Quantifiers

The logical operations indicated by phrases such as “there exists” and “for all” are called *quantifiers* because they indicate the quantity of values for which a proposition is stipulated to hold, not in numerical terms, but in terms of sets. When all of the variables on which the truth value of an open proposition depends are quantified, it becomes a closed proposition. For example, the truth value of the proposition that $x \geq 0$ depends on the value of x , but the proposition that, for all natural numbers x , $x \geq 0$, is true, as is the proposition that there is a real number x such that $x \geq 0$, whereas the proposition that for all real numbers x , $x \geq 0$, is false. Thus, quantification of a variable is a logical operation that removes the dependence of truth value on the specific value of that variable. Instead, the truth value depends on the manner in which the variable is quantified. Quantification creates a proposition with some degree of generality. Since the ability to make general claims lies at the heart of mathematics, understanding quantifiers is absolutely crucial for any mathematics student.

Although, as with other logical operations, phrasing may vary somewhat when propositions are expressed in ordinary language, there are only two types of quantifiers: existential (represented by “there exists”) and universal (represented by “for all”). They may be denoted by the symbols \exists and \forall , respectively. An existential quantifier is always followed by a property satisfied by the object that is proposed to exist; this property is introduced by a phrase such as “such that,” or “for which.” An existential quantifier may also be further specified by a set to which the proposed object must belong; this additional specification just amounts to conjoining an additional property that the object must satisfy. A universal quantifier may be unrestricted, but is often followed by a specified set to which the quantifier applies. The application of a universal quantifier may also be further restricted to elements of that set satisfying a particular property, which just amounts to separating the elements satisfying that property as the set to which the quantifier applies. Similarly, a condition is imposed in a universally quantified proposition amounts to separating elements satisfying the condition as the set to which the quantifier applies. Thus, the following (true) propositions are synonymous:

- For all natural numbers x , $x \geq 0$.
- $\forall x \in \mathbb{N}, x \geq 0$.
 Parentheses are sometimes used instead of a comma: $(\forall x \in \mathbb{N})(x \geq 0)$.

And they are logically equivalent to:

- $\forall x$, if $x \in \mathbb{N}$, then $x \geq 0$.

The following (true) propositions are synonymous:

- For all real numbers x such that $x \geq 0$, $x^2 - x - 2 = 0 \Leftrightarrow x = \sqrt{x+2}$.
- $\forall x \in \mathbb{R}$ such that $x \geq 0$, $x^2 - x - 2 = 0 \Leftrightarrow x = \sqrt{x+2}$.

And they are logically equivalent to:

- $(\forall x \in \mathbb{R})(x \geq 0 \Rightarrow (x^2 - x - 2 = 0 \Leftrightarrow x = \sqrt{x+2}))$. Note that the former versions might be preferable, however, since they don't require delimiters (such as parentheses) to avoid ambiguity. In addition, the former versions are more natural, since the general background condition $x \geq 0$ is set forth in what amounts to a sort of preamble, where it really belongs.

The following (true) propositions are synonymous:

- There exists a real number x such that $x \geq 0$.
- $\exists x \in \mathbb{R}, x \geq 0$. Parentheses or the phrase "such that" may replace the comma.

And they are logically equivalent to:

- $(\exists x)(x \in \mathbb{R} \text{ and } x \geq 0)$.

Quantifiers may be combined with other logical operations; in particular, we must often consider their negations. For example, the negation of the (false) proposition that for all real numbers x , $x \geq 0$, is the (true) proposition that there exists a real number x such that $x < 0$. The negation of the (true) proposition that there exists a real number x such that $x \geq 0$ is the (false) proposition that for all real numbers x , $x < 0$. In general, the proposition that all elements of a certain set have a certain property is negated by the assertion that some element of that set fails to have it; that is, some element of the specified set satisfies the negation of the specified property. The proposition that

some element of a certain set has a certain property is negated by the assertion that all elements of that set fail to have it; that is, all elements of the specified set satisfy the negation of the specified property.

Negating a proposition with multiple quantifiers is a bit like peeling an onion: you have to work your way from the outside in. This procedure may take a bit of practice!

Example. What is the negation of the (true) proposition that for every real number y , there is a real number x such that $x < y$?

Solution. The original proposition asserts a property universal to all real numbers y ; therefore, its negation must assert the existence of a real number y that is an exception. The universal property asserted by the original proposition is that a real number x exists with a particular property relative to the given number y ; therefore, if the real number y is an exception, no number x having this particular property relative to y must exist. To say that no real number x with a particular property exists is to say that all real numbers x fail to have it. The particular property in question is that $x < y$; therefore, to say that x fails to have this property is simply to say that $x \not< y$. Thus the negation of the proposition that for all real numbers y , there is a real number x such that $x < y$, is the (false) proposition that there exists a real number y such that for every real number x , $x \not< y$. The pattern is clearly evident when the statements are written symbolically:

The negation of the statement

$$(\forall y \in \mathbb{R})(\exists x \in \mathbb{R})(x < y)$$

is the statement

$$(\exists y \in \mathbb{R})(\forall x \in \mathbb{R})(x \not< y).$$

□

Remark. In Chapter 5 we will discuss the trichotomy property of order in the real number system, from which it follows that $x \not< y \Leftrightarrow x \geq y$. However, the assumption of trichotomy for this collection of relationships among real numbers lies outside the nature of logic: it is by no means logically necessary for a collection of relationships between elements of a set to satisfy trichotomy. For example, the divisibility relationships among natural numbers do not: 2 does not divide 5, 2 is obviously not equal to 5, and 5 does not divide 2, either.

In general, suppose $p(x, y)$ denotes an open proposition whose truth value depends only on the values of x and y , and $\neg p(x, y)$ denotes the negation of $p(x, y)$. The negation of the statement

$$(\forall x \in A)(\exists y \in B)(p(x, y))$$

is the statement

$$(\exists x \in A)(\forall y \in B)(\neg p(x, y)).$$

Remark. This process may be continued indefinitely with more quantifiers, but statements with a large number of quantifiers get cumbersome and should be formulated in stages, giving names to properties that are defined by intermediate steps.

4.4.2 Relationships Defined by Quantified Propositions

Consider two arbitrary sets A and B . We say that A is a *subset* of B , denoted $A \subseteq B$, if every element of A is also an element of B . Formally:

Definition. $A \subseteq B$ if and only if $(\forall x)(x \in A \Rightarrow x \in B)$. (Equivalently, $(\forall x \in A)(x \in B)$.)

It seems reasonable that the set of all subsets of a given set A is well defined; it is called the *power* set of A and denoted $\mathcal{P}(A)$. The proposition that $x \in \mathcal{P}(A)$ is, by definition, synonymous with the proposition that $x \subseteq A$. The existence of the power set of any given set is asserted by the following axiom:

Axiom (Existence of Power Sets). *Given any set A , $\mathcal{P}(A) = \{x : x \subseteq A\}$ is a set.*

Two sets are deemed to be the same if they have the same elements. Thus, equality between sets is formally specified by the following axiom:

Axiom (Axiom of Extensionality). *For any sets A and B , if $(\forall x)(x \in A \Leftrightarrow x \in B)$, then $A = B$.*

The name of this axiom comes from the notion that sets are determined by the “extent” of what they contain. One might reasonably argue that this is a definition. That it is an axiom reflects the philosophical perspective that we cannot define what it means for things to be equal, but only propose what properties sufficiently characterize them to determine if they are the same. (For example, mountain lions maintain exclusive territories, keeping all other lions out; therefore, if two lions are seen in the same place and no lions have died, they must be the same.) The Axiom of Extensionality declares that sets are determined by their elements. Note that we need not include the converse implication, since it is in the nature of logical expression that if $A = B$, then $(\forall x)(x \in A \Leftrightarrow x \in B)$: if $A = B$, the $x \in A$ and $x \in B$ are simply synonymous. If two things are equal, any proposition about one is synonymous with the proposition obtained by substituting the other.

It is clear from the definition of the subset relationship that the proposition that $A = B$ is synonymous with the proposition that $A \subseteq B$ and $B \subseteq A$. It is sometimes easier to prove two sets are equal by proving each of these conjoined component propositions separately.

The existential proposition that $(\exists x)(x \in A)$ is synonymous with the proposition that $A \neq \emptyset$.

4.5 Exercises

1. Rephrase each implication as a logically equivalent disjunction. Do not use the phrase “it is not the case”; move any negations into the atomic components of your answer.
 - (a) If $x > 2$ then $x^2 > 4$.
 - (b) $x < 2 \Rightarrow x^2 < 4$.
 - (c) $x^2 < 4 \Rightarrow x < 2$
 - (d) If $-2 < x < 2$, then $x^2 < 4$.
2. Write the contrapositive of each proposition in the previous problem.
3. For each of the propositions in the previous problem:
 - Decide if the proposition is true for all real numbers, true for some, but not all, real numbers, or false for all real numbers.
 - For those propositions that are true for some, but not all, real numbers, give an example of a number for which the proposition is true and an example of a number for which it is false.
 - Write the negation of the proposition. Do not use the phrase “it is not the case”; move any negations into the atomic components of your answer.
4. Rephrase each disjunction as a logically equivalent implication. Do not use the phrase “it is not the case”; move any negations into the atomic components of your answer.
 - (a) $x < 2$ or $x > 3$
 - (b) $x < 2$ or $x^2 > 4$

(c) $x > 2$ or $x^2 > 4$

(d) $x > 2$ or $x^2 < 4$

5. There are two fundamental choices for rephrasing a disjunction as a logically equivalent implication, as in the previous problem, depending on which component one negates. What do you notice about them?

6. For each of the propositions in the previous problem:

- Decide if the proposition is true for all real numbers, true for some, but not all, real numbers, or false for all real numbers.
- For those propositions that are true for some, but not all, real numbers, give an example of a number for which the proposition is true and an example of a number for which it is false.
- Write the negation of the proposition. Do not use the phrase “it is not the case”; move any negations into the atomic components of your answer.

7. Write the negation of each of the following propositions without using the phrase “it is not the case”; move any negations into the atomic components of your answer.

(a) $2 < x < 3$

(b) $x > 0$ and $x^2 = 4$

8. Write the negation of each of the following propositions without using the phrase “it is not the case”; move any negations into the atomic components of your answer.

(a) $x > 0$, and $x < 2$ or $x > 5$.

(b) $x > 0$ and $x < 2$, or $x > 5$.

(c) $2 < x < 3 \Rightarrow 4 < x^2 < 9$.

(d) If $x < 2$ or $x > 4$, then $x^2 \neq 9$.

9. For each of the propositions in the previous problem:

- Decide if the proposition is true for all real numbers, true for some, but not all, real numbers, or false for all real numbers.

- For those propositions that are true for some, but not all, real numbers, give an example of a number for which the proposition is true and an example of a number for which it is false.
10. Write the negation of each of the following propositions without using the phrase “it is not the case”; move any negations into the atomic components of your answer.
- (a) $x > 0 \Rightarrow (x^2 = x + 2 \Rightarrow x = \sqrt{x + 2})$.
- (b) $(x > 0 \Rightarrow x^2 = x + 2) \Rightarrow x = \sqrt{x + 2}$.
11. For each of the propositions in the previous problem:
- Decide if the proposition is true for all real numbers, true for some, but not all, real numbers, or false for all real numbers.
 - For those propositions that are true for some, but not all, real numbers, give an example of a number for which the proposition is true and an example of a number for which it is false.
12. For each proposition below:
- Decide if it is true or false.
 - Write its negation without using the phrase “it is not the case”; move any negations into the atomic components of your answer.
- (a) Every real number is greater than 3 or less than 4.
- (b) Every real number is greater than 3 and less than 4.
- (c) Some real number is greater than 3 or less than 4.
- (d) Some real number is greater than 3 and less than 4.
13. For each proposition below:
- Decide if it is true or false.
 - Write its negation without using the phrase “it is not the case”; move any negations into the atomic components of your answer.
- (a) For every real number x , there is a real number y such that $y < x$.

- (b) For some real number x , there is a real number y such that $y < x$.
- (c) For every positive integer x , there is a positive integer y such that $y < x$.
- (d) For some positive integer x , there is a positive integer y such that $y < x$.

14. For each proposition below:

- Decide if it is true or false.
 - Write its negation without using the phrase “it is not the case”; move any negations into the atomic components of your answer.
- (a) For every positive real number x , there exists a real number y such that, for any real number z , $z < y \Rightarrow 2z < x$.
 - (b) For every positive real number x , there exists an integer y such that, for any real number z , $z < y \Rightarrow 2z < x$.
 - (c) For every positive integer x , there exists a positive integer y such that, for any real number z , $z < y \Rightarrow 2z < x$.
 - (d) There exists a real number y such that, for every real number x and every real number z , $z < y \Rightarrow z < x$.

4.6 The Structure of a Proof

Now that we have some precise axioms and definitions, we can write some genuine, complete, rigorous proofs. In a complete proof, every assertion must be justified: if it has not been assumed to be true, we must show that it is logically forced to be true by those propositions that have either been assumed to be true (the axioms) or previously justified. We may use the meaning of logical operations to replace propositions with synonymous propositions, and our definitions allow us to make additional replacements of propositions with synonymous propositions. That’s it! Period! Do not assert anything else in a proof, no matter how obvious it seems! It may be true, and it may be obvious, but you must *explicitly* demonstrate that it logically follows from assertions that have been already established. (Otherwise you have just written a sketch of the proof at best, and worse yet, you may have asserted something false, in which case you don’t even have a correct sketch. Lot’s of seemingly obvious things turn out to be false!)

As you might imagine, there are only limited methods by which the truth of a proposition may be established from that of other propositions. Here is a systematic summary of them, with examples. Of course, most proofs combine several methods and are longer than the examples given here.

4.6.1 From Universal to Specific

If a universal proposition has been established, then it is true in every instance. In the following example, the Axiom of Pairing, which applies universally to any sets A and B , is applied to the instance that $A = B = \emptyset$.

Proposition. $\{\emptyset\}$ is a set.

Proof. By the Existence of the Empty Set (our first axiom), \emptyset is a set. By the Axiom of Pairing, $\{\emptyset, \emptyset\} = \{\emptyset\}$ is a set. \square

4.6.2 Generality through Arbitrariness

In proving a universal proposition, it is usually clearest to write the proof in terms of an arbitrary instance of each variable. Usually the proposition we have to prove is conditional as well as universal: we must show that in every instance satisfying the condition (which may be the conjunction of several component conditions), the consequence is true. Since there is nothing to prove in instances that fail to satisfy the condition, we may assume for the sake of argument that our arbitrary instance does satisfy it. (This assumption is, of course, provisional, as indicated by the phrase “for the sake of argument”; there is certainly no reason to presume that the condition is always true!) In other words, we consider a representative element of the set defined by the condition. (The practice of considering a representative element is somewhat analogous to the journalistic practice of using representative individual stories to convey a general situation, in order to make a stronger emotional connection with the audience - but perhaps this is too much of an interdisciplinary stretch!) Here is an example illustrating how this process may take several steps as one works through the levels of universal quantification.

Proposition. For all sets A and B , $A \setminus B \subseteq A$.

Proof. Since the proposition is universal to all sets A and B , let A and B be arbitrary sets. By definition, to prove that $A \setminus B \subseteq A$, we must prove that every element of $A \setminus B$ is an element of A . So let $x \in A \setminus B$. By definition of $A \setminus B$, $x \in A$ and $x \notin B$. In particular, $x \in A$. \square

A similar method may be used to write an argument based on an existential proposition. If it has been established or assumed that something with a particular property exists, we may consider such an object, represented by a variable, and deduce other propositions from its existence. An example of this method appears in Section 4.6.5, on proof by contradiction, below.

4.6.3 Vacuity

Anything is universally true for the elements of the empty set, since there aren't any. An equivalent way to look at this is that if the condition of an implication is universally false, then the implication is universally true. There is nothing to prove, since the consequence is only required to be true for instances that satisfy the condition. We say the implication is *vacuously* true.

Here are two examples:

Proposition. *There is only one set having no elements.*

Proof. Let A and B be sets with no elements. Then the proposition that $(\forall x)(x \in A \Rightarrow x \in B)$ is vacuously true. Conversely, the proposition that $(\forall x)(x \in B \Rightarrow x \in A)$ is also vacuously true. Thus, by definition, $A = B$. \square

This proposition retroactively justifies our assigning a special symbol to the empty set and using the article “the” when referring to it. Our axiom only asserts the existence of “a” set with no elements, so technically we should have designated such a set with a variable until we proved that only one set had this property; uniqueness is what turns a variable into a constant. However, sticking rigidly to logical order in situations like this is often cumbersome and rather pedantic, so retroactive justifications of this kind are common.

The proof of this proposition illustrates the general method for showing that an object with a given property is unique: consider two objects with this property and prove they are equal.

Proposition. *The empty set is a subset of any set; that is, for any set A , $\emptyset \subseteq A$.*

Proof. For any set A , the proposition that $(\forall x)(x \in \emptyset \Rightarrow x \in A)$ is vacuously true. Thus, by definition, $\emptyset \subseteq A$. \square

Corollary. *For any set A , $\emptyset \in \mathcal{P}(A)$.*

A corollary is a proposition that follows immediately from an established result. This one is an immediate consequence of the definition of the power set; the proof is left to you.

4.6.4 Division into cases

If we know that one of two propositions must be true, we may devise separate arguments in each case. In the following example, the definition of union provides a natural

division into cases. This example also gives additional illustrations of generality through arbitrariness.

Proposition. *For any sets A , B , and C , if $A \subseteq C$ and $B \subseteq C$, then $A \cup B \subseteq C$.*

Proof. Let A , B , and C be arbitrary sets such that $A \subseteq C$ and $B \subseteq C$, and let $x \in A \cup B$. By definition of union, $x \in A$ or $x \in B$.

Case 1: $x \in A$. Since $A \subseteq C$, if $x \in A$ then $x \in C$ by definition of subset; hence, $x \in C$.

Case 2: $x \in B$. By a similar argument, $x \in C$ in this case as well.

In either case, $x \in C$; hence, by definition, $A \cup B \subseteq C$. □

It is certainly possible to divide into more than two cases. However one divides into cases, the key point is to make sure that at least one of the cases must hold for all objects under consideration. A good way to do this is to use a proposition and its negation for each division, as we explain now.

We may assert propositions that are logically *tautological*, that is, forced by their logical construction to be true; for example, the proposition that, for any x , $x = \emptyset$ or $x \neq \emptyset$ is tautologically true, as is the proposition that, for any x , it is not the case that *both* $x = \emptyset$ and $x \neq \emptyset$. Given an open proposition and its negation, one must be true, and the other false, in any instance. Being universally tautologically true, the disjunction of any proposition and its negation provides a foolproof division into cases.

The following example demonstrates this method, along with several other notable practices. In particular, the proposition to be proven is biconditional, but writing every step biconditionally is difficult in this case, so we will prove each component separately. To organize our proof in a readable fashion, we will state each component proposition as a *claim*. That is, we will claim it holds, and then provide a proof that it does.

In general, it is good practice to organize a proof strategically into the major assertions that, taken together, lead to the result. This provides a framework for the ideas and logic, so that the argument is not just a formless sequence of steps (however logically correct it might be). Each major assertion is proposed as a claim, telling the reader (and the author) where you are going with the argument that follows, which usually takes multiple steps. If there are many claims, and especially if there are claims within claims, it helps to number them in outline form. Do not confuse claims with cases: a claim is a proposition you are asserting; a case is a possibility, one of which must be true, although you have no way of knowing which one. Finally, notice that in a long proof the quantifiers are often implied in the exposition rather than repeated over and over.

And so, without further ado, here is an example of a more complicated proof:

Proposition. *The power set of $\{\emptyset\}$ is $\{\emptyset, \{\emptyset\}\}$.*

Proof. By definition, to prove that $\mathcal{P}(\{\emptyset\}) = \{\emptyset, \{\emptyset\}\}$, we must prove that, for all sets x , $x \in \mathcal{P}(\{\emptyset\}) \Leftrightarrow x \in \{\emptyset, \{\emptyset\}\}$.

Claim. $x \in \mathcal{P}(\{\emptyset\}) \Leftrightarrow x \in \{\emptyset, \{\emptyset\}\}$.

Let x be a set, and assume $x \in \{\emptyset, \{\emptyset\}\}$; then $x = \emptyset$ or $x = \{\emptyset\}$.

Case 1: $x = \emptyset$. By the corollary to a previous proposition, $\emptyset \in \mathcal{P}(\{\emptyset, \{\emptyset\}\})$.

Case 2: $x = \{\emptyset\}$. By definition of $\mathcal{P}(\{\emptyset\})$, to prove that $x \in \mathcal{P}(\{\emptyset\})$, we must prove that $x \subseteq \{\emptyset\}$. By definition of subset, to prove that $x \subseteq \{\emptyset\}$, we must prove that, $\forall y, y \in x \Rightarrow y \in \{\emptyset\}$. Let $y \in x$; then $y = \emptyset$. $\emptyset \in \{\emptyset\}$, so $x \subseteq \{\emptyset\}$.

Thus, in either case, if $x \in \{\emptyset, \{\emptyset\}\}$, then $x \in \mathcal{P}(\{\emptyset\})$.

Claim. $x \in \mathcal{P}(\{\emptyset\}) \Rightarrow x \in \{\emptyset, \{\emptyset\}\}$.

Let $x \in \mathcal{P}(\{\emptyset\})$. Since $x \in \mathcal{P}(\{\emptyset\})$, $x \subseteq \{\emptyset\}$, and hence, by definition, $(\forall y \in x)(y \in \{\emptyset\})$. Thus, $(\forall y \in x)(y = \emptyset)$.

Case 1: $x = \emptyset$. In this case, $x \in \{\emptyset, \{\emptyset\}\}$.

Case 2: $x \neq \emptyset$. Then there is an element $y \in x$; therefore, since $(\forall y \in x)(y = \emptyset)$, $x = \{\emptyset\} \in \{\emptyset, \{\emptyset\}\}$.

Thus, in either case, if $x \in \mathcal{P}(\{\emptyset\})$, then $x \in \{\emptyset, \{\emptyset\}\}$

□

The fact that a proposition and its negation are mutually exclusive (that is, only one can be true in any instance) can be helpful in proving that exactly one of several options occurs, as in the following example:

Proposition. *Given any sets A and B , every element of $A \cup B$ belongs to one and only one of the following three sets: $A \setminus B$, $B \setminus A$, or $A \cap B$.*

Proof. Let A and B be sets, and let $x \in A \cup B$. By definition, $x \in A$ or $x \in B$.

Case 1: $x \in A$.

Case 1(a): $x \in B$. Then by definition, $x \in A \cap B$.

Case 1(b): $x \notin B$. Then by definition $x \in A \setminus B$.

Case 2: $x \notin A$. Then, since $x \in A$ or $x \in B$, it must be the case that $x \in B$; hence, by definition, $x \in B \setminus A$.

Thus, we know that x belongs to one of the sets $A \cap B$, $A \setminus B$, or $B \setminus A$. Conversely: if $x \in A \cap B$, then by definition $x \in A$ and $x \in B$ (Case 1(a)); if $x \in A \setminus B$, then by definition $x \in A$ and $x \notin B$ (Case 1(b)); and if $x \in B \setminus A$, then by definition $x \in B$ and $x \notin A$ (Case 3). Since these cases are mutually exclusive, x belongs to only one of these sets. \square

4.6.5 Proof by Contradiction

Our final strategy is to provisionally assume the negation of the proposition we want to prove and, from this assumption, deduce a *contradiction*, meaning a statement that has been established to be false. It therefore follows that our proposition must be true, since its negation led directly to a false proposition. There are a variety of names for this method of proof, including *reductio ad absurdum* (Latin for “reduction to the absurd”), “indirect proof” (since we reach the desired result indirectly, by first supposing it were to be false), or just plain “proof by contradiction.”

The name “proof by contradiction” arises from the fact that one obtains an incontrovertibly false statement either by contradicting (that is, asserting the negation of) a statement that has been established or assumed to be true (for example, $(\exists x)(x \neq x)$, since it is in the fundamental nature of logical expression that $(\forall x)(x = x)$), by asserting an instance satisfying the conjunction of two contradictory propositions (for example, $(\exists x)(\exists A)(x \in A \text{ and } x \notin A)$), or occasionally by asserting the conjunction of contradictory closed propositions (for example, $(\forall x)(x \in A)$ and $(\exists x)(x \notin A)$). The latter two types of propositions, whose falsity rests on their logical construction, are simply the negations of tautologies.

When beginning a proof by contradiction, I prefer to use the word “suppose,” sometimes augmented by the phrase “by way of contradiction” or use of the subjunctive mood, when assuming the negation of the proposition to be proven, to emphasize for clarity the fact that this assumption is not only provisional, but is going to be rejected as false (unlike our provisional assumption of any conditions in a direct proof, which need not be the case, but could be). It is always nice to let the reader know (and remind yourself) what is coming down the pike!

A nifty way to look at proof by contradiction is as follows: Let p be the proposition to be proven. Let F represent a contradiction, and let T represent its negation (which, of course, is true). In proof by contradiction, we show that $\neg p \Rightarrow F$, which is synonymous

with its contrapositive, $T \Rightarrow p$. Since the condition T is true, the consequence p must be as well.

As one might expect, proof by contradiction is the method to be used for proving any proposition that is stated in the negative and cannot be rephrased positively. Thus, it is the method to be used for proving a set is empty, as in the following example:

Proposition. For any sets A and B , $(A \cap B) \cap (A \setminus B) = \emptyset$.

Proof. Suppose by way of contradiction that $(A \cap B) \cap (A \setminus B)$ were not empty. Then there would be some element $x \in (A \cap B) \cap (A \setminus B)$. For this element we would have, by definition of intersection, that $x \in A \cap B$ and $x \in A \setminus B$; therefore, again following the definitions, $x \in A$ and $x \in B$ and $x \in A$ (a redundant assertion we would normally omit) and $x \notin B$. Since $x \in B$ and $x \notin B$ is a contradiction, no such x can exist, so we conclude that $(A \cap B) \cap (A \setminus B) = \emptyset$. \square

Definition. Sets A and B are *disjoint* if (and only if) $A \cap B = \emptyset$.

Using this definition, we may rephrase the previous proposition as stating that, for all sets A and B , the sets $A \cap B$ and $A \setminus B$ are disjoint.

Remark. All definitions are biconditional; sometimes the phrase “and only if” is omitted, but it is always implicit in a definition. This convention is the one exception to the rule in mathematics that all assertions must be explicitly stated.

Students tend to love proof by contradiction, so I hasten to point out that it should not be overdone. Often a direct proof is simpler and clearer. Suppose, for example, that you set out to prove a statement of the form $(\forall x \in A)(p(x) \Rightarrow q(x))$. You suppose by way of contradiction that $(\exists x \in A)(p(x) \text{ and } \neg q(x))$. You consider such an element x and, after a series of deductions in which, as it happens, you never need the component $p(x)$ of your supposition, you deduce $\neg p(x)$, a contradiction. In this situation, it would be easier and clearer to prove directly the contrapositive proposition that $(\forall x \in A)(\neg q(x) \Rightarrow \neg p(x))$. (Presumably, if you didn't use $p(x)$ in your proof by contradiction, you did use $\neg q(x)$, since it is the only other proposition you supposed. If you didn't, then you can prove - directly - the stronger, unconditional statement that $(\forall x)p(x)$.)

As a bad example to demonstrate this, consider a proof by contradiction of our earlier proposition that if $A \subseteq C$ and $B \subseteq C$, then $A \cup B \subseteq C$. Suppose not; that is, suppose that $A \subseteq C$ and $B \subseteq C$ and $A \cup B \not\subseteq C$. Since $A \cup B \not\subseteq C$, there is an element $x \in A \cup B$ such that $x \notin C$. Since $x \in A \cup B$, $x \in A$ or $x \in B$. In the first case, $A \not\subseteq C$, contradicting a component of our hypothesis. Similarly, in the second case, $B \not\subseteq C$, again contradicting a component of our hypothesis. We conclude that if $A \subseteq C$ and $B \subseteq C$, then $A \cup B \subseteq C$. Even though we have abbreviated our exposition in comparison to the

laborious style of the earlier proof, in which we strove to make every logical step explicit, the awkwardness of this proof should nonetheless be apparent. Hypotheses that are not used directly are confusing, hence bad style, and too many negatives give anyone over the age of two a headache!

4.7 Exercises

In the following exercises, you may use any of the propositions proven in the examples.

1. Prove that $\{\emptyset, \{\emptyset\}\}$ is a set.
2. Prove that $\{\emptyset, \{\emptyset\}, \{\emptyset, \{\emptyset\}\}\}$ is a set.
3. Prove that, for any sets A and B , $A \setminus B$ is indeed a set. (Hint: Use the Axiom of Restricted Comprehension.)
4. Prove that, for any sets A and B , $A \cap B$ is a set.
5. Prove that, for any sets A and B , $A \cap B = A \setminus (A \setminus B)$.
6. Prove that for any set A , $A \subseteq A$. This result explains why we prefer to use the symbol \subseteq for the subset relationship rather than \subset .
7. Prove that for any set A , $\emptyset \in \mathcal{P}(A)$.
8. Prove that for any sets A and B , $A \cup B = B \Leftrightarrow A \subseteq B$.
9. Prove that for any sets A and B , $A \setminus B$ and $B \setminus A$ are disjoint.
10. Prove that for any sets A and B , $A \subseteq A \cup B$.
11. Prove that for any sets A and B , $A \cap B \subseteq A$.
12. Prove that for any sets A and B , $A = (A \setminus B) \cup (A \cap B)$.
13. Prove that for any sets A and B , $(A \cup B) \setminus B = A \setminus B$.
14. Prove that for any sets A and B , $A \setminus (A \cap B) = A \setminus B$.
15. Let A , B , and U be any sets such that $A \subseteq U$ and $B \subseteq U$. (The letter “ U ” is used because we can think of U as the “universe” to which the elements of sets A and B belong.) The following propositions are called De Morgan’s Laws. Prove them.

- (a) $A \cap B = U \setminus ((U \setminus A) \cup (U \setminus B))$.
- (b) $A \cup B = U \setminus ((U \setminus A) \cap (U \setminus B))$.
16. Prove that $A \cap B = (A \cup B) \setminus (A \ominus B)$. (Hint: Apply one of De Morgan's laws, choosing an appropriate set U and using previous results to establish the necessary hypothesis. This approach results in a much shorter proof than starting from scratch. Use previous results whenever possible!) The current result shows that the intersection operation can be defined in terms of union and complementation.
17. Write $A \setminus B$ using only the operations of union (\cup) and symmetric difference (\ominus). Prove your expression is correct.
18. Is it possible to write $A \ominus B$ using only the operations of union (\cup) and intersection (\cap)? If so, do it and prove your expression is correct. If not, explain why not.
19. Consider the operations \cup , \cap , \ominus , and \setminus as binary operations on sets.
- (a) Which are associative? Prove your answers. (That is, for those that are, prove the appropriate universal proposition, and for those that aren't, provide a specific counterexample. A counterexample is simply an instance that disproves a universal proposition by showing the existence of an exception.)
- (b) Which are commutative? Prove your answers.
- (c) Which pairs of operations satisfy a distributive relationship. Prove your answers.
20. Compute the power set of $\{\emptyset, \{\emptyset\}, \{\emptyset, \{\emptyset\}\}$. For **extra credit**, prove your answer is correct.
21. Now that we have added an axiom asserting the universal existence of power sets and the Axiom of Extensionality to our collection, can we prove the existence of an infinite set? Explain.

Chapter 5

The Real & Natural Number Systems

In which we logically develop the properties of the number systems used to measure continuous and discrete quantities.

Set theory and logic establish the basis for mathematical thought. Mathematics beyond these realms generally concerns the elements of some set or sets with particular properties, for example, real numbers (as in algebra), natural numbers (as in number theory), or points or lines in some geometry (whether planar, spherical, on some other surface, or in a higher dimension). As we will see, the relations among the elements of these sets, such as the association of a sum and product to each pair of real numbers or the incidence relation between points and lines, which establishes which lines pass through which points, are also given by sets with particular properties. The properties the sets in question are assumed to have, called axioms, are essential to proving results about them. In order to write proper, rigorous proofs, we must know and state precisely and explicitly what these properties are. Many properties, such as associativity in various algebraic systems, recur in different contexts, so we name them in order to systematize our knowledge and better see parallels and connections. It is also important that the basis for our axioms be examined, so that it doesn't seem as if we are just playing games with rules made up out of thin air.

We will begin with careful study of the real and natural number systems, which are not only central to mathematics, but also to science and everyday life. Real numbers are used to measure continuous quantities such as length, area, volume, and mass. Natural numbers are used to count discrete things. (A note on English usage: The word “fewer”, not the word “less,” is used in comparing discrete quantities, as in, “I have fewer cookies

than you have.” The word “less” is reserved for continuous quantities, as in, “My house has less space than yours.” Using these words properly helps emphasize and clarify what kind of numbers we are talking about.) You have undoubtedly had much experience solving algebra problems involving real numbers, yet you have probably not proven many of the results you applied. You also regularly use many basic facts about the natural numbers, such as factorization into primes, when doing arithmetic. (Think about reducing fractions, for example) Yet you probably have not thought hard about what general properties you are using, why they are true, or in what ways they extend or fail to extend to other elementary systems with operations, such as polynomials. In your future classes, you will have opportunities to study the theories of geometry, advanced algebra, and perhaps advanced number theory, but the basic algebra and arithmetic with which you are familiar will likely be assumed throughout the remainder of your education without discussion. One day not too long from now, many of you will teach these basic subjects. This course may be the only opportunity you have to thoroughly study the foundation on which they are built.

The natural numbers are, as we all know, a subset of the real numbers. On the one hand, this makes them more limited; certain operations, such as division and subtraction, cannot always be carried out. On the other hand, the fact that the natural numbers are a smaller (but still infinite) set with a very special structure makes it possible to define some interesting concepts, such as primeness, that do not make sense for real numbers in general (since every real number is divisible by every other real number except zero) and to prove some fascinating assertions about them. For historical reasons, the theory of the natural number system is called, simply, *number theory*. Many theorems in number theory are easily stated but surprising. Some simply-stated conjectures have only recently been proven, and others remain unproven to this day, centuries after they were first formulated. The existence of such intricate and elusive patterns in a number system that arises simply from counting is a mystery that has fascinated people since the earliest times.

In developing these number systems, we have two options. One is to start with the natural numbers and build, first the integers, then the rational numbers, and then the real numbers out of them. The other is to start with the real number system as a whole and define the natural numbers as a specific subset. (The integers and rational numbers can then be defined as subsets of \mathbb{R} based on the natural numbers: the integers consist of the natural numbers and their additive inverses; the rational numbers consists of those real numbers equal to the ratio of two integers. Formally, $\mathbb{Z} = \mathbb{N} \cup \{-n : n \in \mathbb{N}\}$, and $\mathbb{Q} = \{\frac{m}{n} : m \in \mathbb{Z} \text{ and } n \in \mathbb{Z}\}$. We will follow the convention that $0 \in \mathbb{N}$; hence 0 need not be included separately in the set of integers.) Humans seem to have a built in natural intuition, much like our language ability, about both continuous quantities,

such as length and area, and discrete quantities, such as numbers of objects, and it is not clear that either approach is better. Both approaches have their virtues, and we will in fact discuss both of them. A major virtue of starting with the real number system as a whole is that from a logical point of view this approach is faster, simpler and easier. So we'll take it first.

5.1 Addition and Multiplication of Real Numbers

The set of real numbers, \mathbb{R} , comes with two *binary operations*, addition and multiplication. It is these operations that make \mathbb{R} into a number *system* (as opposed to just a set). We will not define these operations; we will simply assume they exist and obey certain axioms. (Later we will define them in terms of set operations.) These operations are abstractions of our experience working with and thinking about continuous quantities, and it is this experience that provides the basis for the axioms we will assume about them.

Addition reflects our experience with combining like quantities. If we place a 7 foot board end-to-end with a 3 foot board, they extend a length of 10 feet. If we saw a 3 foot board from a 10 foot board with a saw that creates a kerf .01 feet wide, we are left with a 6.99 foot board. It requires $2 + \pi$ meters of string to border a semicircle of radius one meter. If we pour ten liters of water into a tub and then remove the amount needed to fill a hemispherical bowl of radius 10 centimeters, we have $10 - \frac{4\pi}{3}$ liters left in the tub. When we first learn addition, we generally think of combining discrete quantities, as when we combine seven chocolate chip cookies with three oatmeal cookies to obtain ten cookies. This interpretation makes sense for integers and, more generally, for integer multiples of any common unit, but not for sums such as $\sqrt{2} + \sqrt{3}$. (The numbers $\sqrt{2}$ and $\sqrt{3}$ are called *incommensurable*: no matter how each number is divided into equal parts, the parts have to be a different size, so there is no common unit with which they can be measured. If there were, $\sqrt{\frac{2}{3}}$ would be rational, which it demonstrably is not. See if you can prove that!) Integers represent whole units of a quantity, but continuous quantities such as distance or mass occur in all amounts in between. Exactly what is meant by “all amounts in between” will be discussed in Chapter ???. To understand the phrase “in between,” we need to understand order, and to understand “all amounts,” we need to understand a special property of order in the real number system called the least upper bound property. (The fact that the square roots of any rational number are in the real number system is due to this property.)

Multiplication reflects our experience with several independent dimensions. If we use a 15 kilowatt bulb for three hours, we have consumed 45 kilowatt-hours of power. If we

travel at 100 kilometers per hour for 3 hours, we have traveled 300 kilometers. Near the surface of the earth, objects under the influence of gravity accelerate at approximately -9.8 meters per second per second, that is, -9.8 meters/second². (The acceleration is negative on account of the convention that the positive direction is up.) If an object starts at rest and falls for 2 seconds, it will be traveling at -19.6 meter/second. The force required to lift an object with a mass of 2 kilograms is 19.6 kilogram-meters/second²; the unit of a kilogram-meter/second² is called a newton. If you lift the object a distance of 3 meters, you have done 58.8 newton-meters of work; if you lower it by 3 meters, you have done -58.8 newton-meters of work. By exerting a negatively directed force on an object moving in a negative direction, gravity did positive work. If a rectangular field is 10 meters long and $2\sqrt{3}$ meters wide, it has an area of $20\sqrt{3}$ meters². (Note: Insistence on right-angled corners - that is, *square* corners - when multiplying linear dimensions is merely a convenient convention; any angle between the dimensions would do as long as we did not vary it, but our units of area would no longer be *square* units.) When we first learn to multiply, we generally think of repeated addition of a discrete quantity, as when we combine three bags of cookies containing five cookies each to obtain 15 cookies. This interpretation only makes sense for integers, but we can think of the quantities of bags and cookies as representing independent dimensions, obtaining a perspective that is consistent with our continuous general interpretation: 3 bags at 5 cookies per bag yields 15 cookies.

A binary operation associates to each pair of elements of a set, real numbers in this case, an element of the same set, often referred to as the result of the operation. In the case of addition, this result is called the *sum* of the pair; in the case of multiplication, it is called their *product*. Every pair of real numbers has a sum and a product. Given two real numbers a and b , their sum is denoted by $a + b$ and their product by $a \cdot b$ or, when no confusion will result, simply by ab . Because the sum and product are determined uniquely by the numbers themselves, the variables used to represent the numbers do not affect the result: if $a = b$ and $c = d$, then $a + b = c + d$ and $ab = cd$. It is not necessary to give written justification in a proof when making a substitution of this type.

The qualities above constitute the general definition, in non-technical language, of a binary operation. Since that term captures these qualities, we won't list them as separate axioms. We now list the additional properties of addition and multiplication of real numbers.

Axioms of Addition and Multiplication in the Real Number System

There exist binary operations $+$ and \cdot on \mathbb{R} , such that:

1. These operations are *associative*:

(a) $\forall x, y, z \in \mathbb{R}, (x + y) + z = x + (y + z).$

(b) $\forall x, y, z \in \mathbb{R}, (x \cdot y) \cdot z = x \cdot (y \cdot z).$

2. These operations are *commutative*:

(a) $\forall x, y \in \mathbb{R}, x + y = y + x.$

(b) $\forall x, y \in \mathbb{R}, x \cdot y = y \cdot x.$

3. Each operation has a distinct *identity element*:

(a) There exists an element $0 \in \mathbb{R}$ such that, $\forall x \in \mathbb{R}, x + 0 = x.$

(b) There exists an element $1 \in \mathbb{R}$ such that $1 \neq 0$ and, $\forall x \in \mathbb{R}, x \cdot 1 = x.$

4. All possible *inverses* exist:

(a) For each x in \mathbb{R} , there exists a y in \mathbb{R} such that $x + y = 0.$

(b) For each x in \mathbb{R} different from 0, there exists a y in \mathbb{R} such that $x \cdot y = 1.$

5. The operation \cdot *distributes* over $+$:

$\forall x, y, z \in \mathbb{R}, x \cdot (y + z) = (x \cdot y) + (x \cdot z).$

These axioms are motivated by our physical experience of numbers. The basis for each can be readily illustrated. For example, the Figure 5.1 illustrates the basis for the distributive property (Property 5).

5.2 Exercises

1. Prove there is the unique element $y \in \mathbb{R}$ such that $\forall x \in \mathbb{R}, x + y = x.$ (Property 3a establishes that such an element exists. Suppose $y' \in \mathbb{R}$ also satisfies the property that $\forall x \in \mathbb{R}, x + y' = x.$ Consider $y + y'$ to prove that $y = y'.$)

The uniqueness of the element with this property justifies our designating it by the special symbol 0; there is no ambiguity regarding to which element 0 refers.

2. Similarly, prove there is a unique element of $y \in \mathbb{R}$ such that $\forall x \in \mathbb{R}, x \cdot y = x.$

The uniqueness of the element with this property justifies our designating it by the special symbol 1.

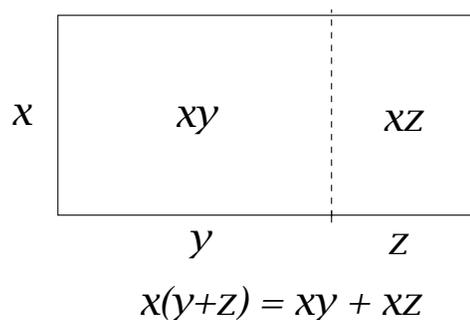


Figure 5.1: Illustration of the basis for the distributive property.

3. Let $x \in \mathbb{R}$ be a fixed number. Prove there is a unique element $y \in \mathbb{R}$ with the property that $x + y = 0$. (Suppose $x + y = 0$ and $x + z = 0$. Show that $y = z$.)

We define the *additive inverse of x* to be this unique number associated to x and denote it by $-x$. We define *subtraction* by $x - y = x + (-y)$.

4. Similarly, given $x \neq 0 \in \mathbb{R}$, prove there is a unique element $y \in \mathbb{R}$ such that $xy = 1$.

We define the *multiplicative inverse of x* to be this unique number associated to x and denote it by $\frac{1}{x}$ or, extending the properties of exponential notation, by x^{-1} . (We will defer further discussion of exponential notation until Section 5.7, on recursive definition.) We define *division* by

$$x \div y = x/y = \frac{x}{y} = x \cdot \frac{1}{y}.$$

Notice the several notations commonly used for the result of division. We will generally use the second or third notation. The last expression in this chain of equalities gives the definition, in terms of the multiplicative inverse and the already

postulated operation of multiplication. Note also that, for any $y \neq 0 \in \mathbb{R}$,

$$1 \div y = 1 \cdot \frac{1}{y},$$

by definition, where $\frac{1}{y}$ means the multiplicative inverse of y , and in turn

$$1 \cdot \frac{1}{y} = \frac{1}{y},$$

by the multiplicative identity property of 1, so our third notation for division of 1 by another number is consistent with the notation we previously chose for the multiplicative inverse. (This consistency is important! If the same symbol could represent two distinct numbers, our meaning would not be clear from the context in many situations.)

5. Prove that for any $x \in \mathbb{R}$, $0 \cdot x = 0$. (Hint: use the distributive property.)
6. Conclude from the preceding result that 0 does not have a multiplicative inverse. (The axioms do not assert this. They guarantee that there is a multiplicative inverse for any number other than zero, but not the converse proposition that there is not one for zero. The proof requires the fact that $1 \neq 0$, asserted as part of Property 3b.)
7. Prove the converse of Exercise 5, concluding that for any real numbers x and y , $xy = 0 \Leftrightarrow (x = 0 \text{ or } y = 0)$. (To construct the converse of Exercise 5, rephrase it as follows: if $x = 0$ or $y = 0$, then $xy = 0$. This statement is slightly stronger than the result of Exercise 5, but is obtained from it simply by applying the commutative property of multiplication. The proposition to be proved in this exercise, then, is that if $xy = 0$, then $x = 0$ or $y = 0$. Hint: Rephrase the conclusion as $x \neq 0 \Rightarrow y = 0$.)
8. Prove that $\{x \in \mathbb{R} : (x - a)(x - b) = 0\} = \{a, b\}$.
9. Prove that for all real numbers x and y , $x = 0 \Leftrightarrow x + y = y$.
10. Prove that for all real numbers x and y , $(x = 1 \text{ or } y = 0) \Leftrightarrow xy = y$.
11. Prove that $(-1)x = -x$. Here we have omitted the obvious implied quantifier ($\forall x \in \mathbb{R}$) for the sake of brevity, as is common practice. (The left side of the equation denotes the product of x and the additive inverse of 1; the right side denotes the additive inverse of x . These are not a priori the same number. That they are requires proof.)

12. Prove that $\frac{x}{y} \cdot \frac{w}{z} = \frac{xw}{yz}$. (Hint: use the commutative property.)
13. Prove that $\frac{x}{y} + \frac{w}{z} = \frac{xz+yw}{yz}$. (Hint: use the distributive property.)

5.3 The Ordering of the Real Numbers

The real numbers also come with an *order relation*. Any two distinct real numbers are comparable, and the result of that comparison is denoted by placing the symbol $<$ between them in the correct order. As with the operations of addition and multiplication, we will not define how to compare numbers, only assume that the order has certain properties motivated by the comparison of physical quantities. (We will later define the order relation in the course of building up the number systems from the natural numbers.) By definition, $y > x$ means $x < y$ and $x \leq y$ means $x < y$ or $x = y$. Also, $x < y < z$ means $x < y$ and $y < z$.

The defining properties of an order relation are:

- Anti-reflexivity: $\forall x \in \mathbb{R}, x \not< x$ (nothing is less than itself); and
- Transitivity: $\forall x, y, z \in \mathbb{R}, x < y$ and $y < z \Rightarrow x < z$.

From anti-reflexivity and transitivity we can deduce a third property satisfied by any order:

- *Anti-symmetry*: $\forall x, y \in \mathbb{R}, x < y \Rightarrow y \not< x$.

Proof. By way of contradiction, suppose there exist numbers x and y such that $x < y$ and $y < x$. By transitivity (one of the assumed properties of order) it follows that $x < x$. But this contradicts anti-reflexivity (the other assumed property of order). We therefore conclude that no such x and y exist; hence, $\forall x, y \in \mathbb{R}, x < y \Rightarrow y \not< x$. \square

Anti-reflexivity and transitivity (and hence anti-symmetry) are properties of any order relation by definition (they go with the territory, so to speak), so we won't list them as separate axioms. However, one can conceive of order relations in which not every two elements of the set being ordered can be compared. A standard example is the ordering of a family of sets by proper inclusion: Let \mathcal{C} be any family of sets (such as the power set of a given set). If $A, B \in \mathcal{C}$, define $A \subset B$ to mean that A is a subset of B but not equal to B . A is said to be a *proper* subset of B . You can readily check that

proper inclusion is anti-reflexive (by definition, since *proper* means that the sets are not equal) and transitive, so it is an order relation. However, it is certainly possible that, given two distinct sets in our family, neither one is a proper subset of the other. Order relations in which some pairs of elements are not comparable are called *partial*. Those in which every element is comparable to every other (*comparability* holds) are called *total*.

To avoid having to refer separately to the algebraic and order properties, we will begin numbering the order properties where we left off.

Axioms of the Order of Real Numbers with Respect to the Operations

There is a total order relation $<$ on the real number system with the following properties:

6. The operations respect the order relation:

(a) If $x > y$, then $x + z > y + z$.

(b) If $x > y$ and $z > 0$, then $x \cdot z > y \cdot z$.

7. The order relation $<$ has the *least upper bound property*.

The least upper bound property is the one that establishes the existence of “every amount in between.” It is needed to show the existence of irrational numbers such as $\sqrt{2}$ and π . All of the integers are obtained from 0 and 1 using only addition and the existence of additive inverses. All of the rational numbers are obtained from the integers using only multiplication and the existence of multiplicative inverses. However, no irrational numbers can be proven to exist using only the algebraic axioms. The least upper bound property is somewhat complicated to define; we will define and discuss it in Chapter ??, after we have thoroughly discussed order relations. It is not needed for algebraic calculations or comparisons, so for now we will not use it.

You are probably familiar with the following property from your high school algebra classes. Now we see that it can be proven from the properties we have assumed.

Theorem (Trichotomy Theorem). $\forall x, y \in \mathbb{R}$, *exactly one of the following is true: $x = y$, $x < y$, or $x > y$.*

Proof. Since $<$ is a total order, we know one of the statements must be true. We need to show that, in each case, the other two statements cannot hold.

Case 1: $x = y$. Then, by anti-reflexivity, $x \not< y$ and $y \not< x$.

Case 2: $x < y$. By anti-symmetry, $x \not> y$. By anti-reflexivity (rephrased contrapositively), $x \neq y$.

Case 3: $x > y$. Similar to Case 2.

□

Note that the Trichotomy Theorem holds for any total order that might be defined on any set; we have not used anything special about \mathbb{R} . (We did not need axiom (6) or (7).)

We have discussed operations and order relations in non-technical language. We will learn shortly that relations can be defined in terms of set theory, and that order relations are just those special relations satisfying anti-reflexivity and anti-transitivity. We will also learn that functions are special relations satisfying a certain property, and that binary operations are special types of functions.

Finally, we can now define what it means for a number to be *positive* or *negative*:

Definition. A real number x is *positive* if $x > 0$; a real number x is *negative* if $x < 0$.

We can summarize the second part of Axiom (6) by saying that multiplication by a positive number preserves order. We will prove (and therefore need not state as an axiom) that multiplication by a negative number reverses order. (Multiplication by zero loses the order information by making both sides equal to zero.)

5.4 Exercises

Assume Axioms (1)-(6) given above and all previously proven results. Prove each of the given propositions.

1. $\forall x, y \in \mathbb{R}, x < y \Leftrightarrow y - x > 0$. (That is, x is less than y if and only if $y - x$ is positive, where positive means greater than zero, as stated in the definition above.)
2. $\forall x \in \mathbb{R}, x > 0 \Leftrightarrow -x < 0$. (That is, a number is positive if and only if its additive inverse is negative. Since $-(-x) = x$, it follows that a number is negative if and only if its additive inverse is positive.)
3. $\forall x, y \in \mathbb{R}, x < y \Leftrightarrow -x > -y$.

Remark. This result incorporates and generalizes the preceding one, since $-0 = 0$. Its proof is just a generalization of the previous proof as well.

4. $\forall x, y, z \in \mathbb{R}, (x < y \text{ and } z < 0) \Rightarrow zx > zy$.
5. (a) The product of two positive numbers is positive.
(b) The product of a negative number with a positive number is negative.
(c) The product of two negative numbers is positive.

We can summarize these results as follows: the product of two numbers is positive if and only if both are positive or both are negative. (We know from a previous result that if one is negative and the other positive, the resulting product cannot be zero, so it must be negative.)

6. $\forall x \in \mathbb{R}, x \neq 0 \Leftrightarrow x^2 > 0$. (In English, the forward implication says that the square of a nonzero real number is positive. For the converse implication, use the contrapositive.)
7. $-1 < 0 < 1$.
8. A number and its multiplicative inverse are either both negative or both positive.
9. (a) $\forall x, y \in \mathbb{R}, 0 < x < y \Leftrightarrow 0 < \frac{1}{y} < \frac{1}{x}$.
(b) $\forall x, y \in \mathbb{R}, 0 > y > x \Leftrightarrow 0 > \frac{1}{x} > \frac{1}{y}$.

(To prove the converse implications, just apply the implication you already proved, using the obvious fact that $\frac{1}{(\frac{1}{x})} = x$. Make sure this fact really is obvious to you!)

10. $\forall w, x, y, z \in \mathbb{R}, (w < y \text{ and } x < z) \Rightarrow w + x < y + z$.

Is the converse true? If so, prove it. If not, give a counterexample.

11. The sum of two positive numbers is positive and the sum of two negative numbers is negative.
12. Between any two real numbers, there is another real number. Formally, $\forall x, z \in \mathbb{R}, \exists y \in \mathbb{R}$ such that $x < y < z$. (Hint: Prove that the obvious candidate, $y = \frac{x+z}{2}$, does the job.)

5.5 The Definition of the Natural Number System

We all know what the natural numbers are: $\{0, 1, 2, 3, \dots\}$. (There are differing opinions on whether or not zero should be included - we include it. ¹) But suppose someone asked you to prove that some proposition was true for all natural numbers. What would you do? As a subset of the real number system, the natural numbers inherit all the axiomatic properties that are true for all real numbers, but the proposition to be proven might not apply to all real numbers; rather, its proof might require something special about the natural numbers.

For example, the following proposition is true: For any natural number n , the sum of all the natural numbers less than n is $\frac{n(n+1)}{2}$. This proposition does not even make sense for all real numbers, because the set of real numbers less than a given number is not finite. (Even an infinite series would not make sense, because the set of real numbers less than a given number cannot be listed, even in an infinite sequence.) How could this proposition be proven? “Let n be a natural number. Then \dots ” what, exactly? We can’t check the sum for every number on the list!

An infinite list suggested by a pattern is not adequately precise for proofs. We need a definition for the set of natural numbers. The properties in this definition can then be used in proofs. To formulate a definition, we must identify the essential - that is, defining - properties that the natural numbers, as we conceive them, possess.

We recognize that the system of natural numbers, as the system used for counting discrete objects, has three essential properties:

- It has a specified starting point, the first and smallest natural number (0 in our case, marking the absence of any object, 1 if we were to take counting to begin with the first object present).
- Adding 1 to any natural number gives a new natural number. No matter how large the crowd, you can always add one more. Once we learn a systematic way of naming the numbers, we could explicitly count forever; even though we might not know an English word for every number, once we decide on a place value system we can write the *numeral* for any number. (The really large ones would take a very, very long - but finite - piece of paper. If you use the binary system, for example, you only need two symbols for enumeration, “0” and “1”, but the numerals quickly get long: 0, 1, 10, 11, 100, 101, 110, 111, 1000, \dots)

¹It is true that zero is somewhat different, because we do not use it to count objects, but rather to indicate that there are no objects to count. Historically, zero arose later than the other natural numbers. Nonetheless, logicians often include zero in the natural number system because, in constructing representatives of the natural numbers using set theory, we begin with the empty set.

- The only natural numbers are the ones that can be reached from 0 (or 1) by counting, adding one at a time. There are no others.

To more succinctly refer to the second property above, we give it a name:

Definition. A set S of real numbers is *inductive* if, $\forall s \in S, s + 1 \in S$. (You might prefer to think of this property in its conditional form as a statement about all real numbers: $\forall s \in \mathbb{R}, s \in S \Rightarrow s + 1 \in S$.)

We can now define the set of natural numbers:

Definition. The set of natural numbers, denoted by \mathbb{N} , is the intersection of all inductive subsets of \mathbb{R} that contain 0.

This definition captures, with precision, exactly the three properties listed above. The intersection of a family of sets, each of which contains 0, contains 0. (The proof which follows directly from the definition of the intersection of a family of sets, is left to the exercises.) The intersection of a family of inductive sets is inductive (a fact whose proof is also left to the exercises). By taking the intersection of all sets of real numbers with these properties, we obtain the smallest one, including only those numbers that are absolutely necessary to obtain an inductive set, starting with 0.

5.6 Proof by Induction

A direct consequence of the definition of the set of natural numbers is that any inductive set that contains 0 must contain the entire set of natural numbers. (See Exercise 4.) Therefore, to show that a proposition is true for all natural numbers, we simply need to show that the set of numbers for which the proposition is true includes 0 and is inductive. This proof strategy is called *proof by induction*. We now demonstrate it with a simple proposition.

Proposition. *Every natural number is greater than or equal to 0.*

Proof. We examine the set $[0, \infty) = \{x \in \mathbb{R} : x \geq 0\}$, showing that it contains 0 and is inductive. (In naming this set, I have used the customary interval notation.)

Claim. $0 \in [0, \infty)$.

This is obvious: $0 = 0$.

Claim. *The set $[0, \infty)$ is inductive.*

This means to show that if $x \in [0, \infty)$, then $x + 1 \in [0, \infty)$. So let $x \in [0, \infty)$. Then, by the defining property of $[0, \infty)$, $x \geq 0$. We proved earlier that $1 > 0$; hence, by Axiom 6 of the real number system and the defining property of 0 , $x + 1 > x + 0 = x$. So by transitivity (or substitution in the case $x = 0$), $x + 1 > 0$. Thus $x + 1 \in [0, \infty)$.

Since $[0, \infty)$ is inductive and contains 0 , $\mathbb{N} \subseteq [0, \infty)$. In other words, every natural number is greater than or equal to 0 . \square

In writing the proof above, I have emphasized by my phrasing that we are proving that the set specified by a certain property contains the set of natural numbers; again, we do this by proving that the specified set contains 0 and is inductive (so it is one of those sets whose intersection creates the natural numbers). It is helpful to remember that all of the variants of proof by induction are just modifications or refinements of this idea.

It is customary to write proofs by induction in an abbreviated form by simply considering the property in question. We may also restrict our attention to natural numbers with that property; thus, we are showing that the set of natural numbers with the specified property contains, and is therefore equal to, the whole set of natural numbers. For the proposition above, it was not necessary to assume this restriction, since the entire set of non-negative real numbers is inductive; however, often the restriction to natural numbers is needed, as the set of all real numbers with the desired property may not be inductive, or even definable. (For example, as noted above, it makes no sense to define the set of real numbers x such that the sum of all of those numbers less than or equal to x is $\frac{x(x+1)}{2}$.) Here is the same proof written in the customary form:

Proof.

Claim. $0 \geq 0$. *Obvious.*

Claim. $n \geq 0 \Rightarrow n + 1 \geq 0$. (*Actually, with no extra work we get the stronger statement that $n \geq 0 \Rightarrow n + 1 > 0$.*)

Assume $n \geq 0$. We proved previously that $1 > 0$. Thus, by Axiom 6 of the real number system and the defining property of 0 , $n + 1 > n \geq 0$. (Note: We could just have well used the alternative argument that $n + 1 \geq 1 > 0$.) \square

Remark. The first claim is customarily called the *initial* claim (or the *initial step* of the proof), and the second claim is customarily called the *inductive* claim (or inductive step). The hypothesis assumed to prove the conclusion of the inductive claim is called the *inductive hypothesis*.

Remark. It is important to recognize that we are *not* making a circular argument by assuming that $n \geq 0$ in the inductive claim. The proposition we are proving, $(\forall n \in \mathbb{N})(n \geq 0)$, is universal and *unconditional*. The inductive claim, $(\forall n \in \mathbb{N})(n \geq 0 \Rightarrow n + 1 \geq 0)$ is also universal, but *conditional*. In assuming as our hypothesis its condition that $n \geq 0$, we are merely considering an arbitrary *particular* (and hypothetical) natural number n that would satisfy the condition, not making an assumption about all natural numbers.

5.7 Recursive Definition

We would like to now explore the proposition that, for any natural number n , the sum of all the natural numbers less than n is $\frac{n(n+1)}{2}$. However, we have a technical problem. We all know what is meant by the sum of all the natural numbers less than n , but the number of terms in this sum depends on n and can be arbitrarily large. How do we define it precisely? (As you know, nothing can be rigorously proven about a concept that is not precisely defined!) What we need is called a *recursive* definition:

Definition.

$$\sum_{i=0}^n i = \begin{cases} 0, & \text{if } n = 0 \\ (\sum_{i=0}^{n-1} i) + n, & \text{if } n > 0. \end{cases}$$

Remark. In other words, we have defined each multiple sum in terms of a binary sum with the previously defined result: $\sum_{i=0}^0 i = 0$, $\sum_{i=0}^1 i = 0 + 1 = 1$, $\sum_{i=0}^2 i = 1 + 2 = 3$, $\sum_{i=0}^3 i = 3 + 3 = 6$, \dots . This agrees, of course, with the way we would add up a string of numbers in practice. The word “recursive” comes from the Latin root *cursare*, meaning *to run*, and prefix *re-*, meaning *back*. We “run back” to get the preceding output in order to use it in obtaining the next one.)

Remark. More informally, we can write $\sum_0^n i = 0 + 1 + 2 + 3 + \dots + n$. Anytime an ellipsis (“ \dots ”) is used to indicate a pattern involving the natural numbers up to some arbitrary point, a recursive definition is really being made. Any time an expression with an ellipsis is used in a calculation, a proof by induction is really being informally conveyed.

Now that we have defined $\sum_{i=0}^n i$, we can prove our proposition:

Proposition. For all natural numbers n , $\sum_{i=0}^n i = \frac{n(n+1)}{2}$.

Proof.

Claim. *The result is true for 0.*

$$\sum_{i=0}^0 i = 0 = \frac{0(1)}{2}$$

by definition and basic arithmetic (results proven previously).

Claim. *If the result is true for a natural number $n - 1$, then it is true for n .*

Assume

$$\sum_{i=0}^{n-1} i = \frac{(n-1)(n)}{2}.$$

By definition,

$$\sum_{i=0}^n i = \left(\sum_{i=0}^{n-1} i \right) + n,$$

and by hypothesis,

$$\sum_{i=0}^{n-1} i = \frac{(n-1)(n)}{2}.$$

Thus, substituting, we obtain

$$\sum_{i=0}^n i = \frac{(n-1)(n)}{2} + n = \frac{n^2 - n + 2n}{2} = \frac{n^2 + n}{2} = \frac{n(n+1)}{2}.$$

□

The two claims in the proof above show that the set of numbers n for which $\sum_{i=0}^n i = \frac{n(n+1)}{2}$ contains 0 and is inductive, respectively. It follows that this set contains all the natural numbers: the equation $\sum_{i=0}^n i = \frac{n(n+1)}{2}$ holds for all of them.

Remark. It is often convenient to use the expression $n - 1$ for the arbitrary natural number in our condition in order to match the notation of a recursive definition, in which case the next number is n . Since $n - 1$ is a natural number, $n \geq 1$. You can just as well use n and $n + 1$ if you prefer (in which case n may be any natural number) as long as you are consistent.

Proof by induction is an incredibly powerful method, so powerful that its power can be a drawback to understanding. The logic proceeds as if by magic, sometimes giving little intuition into *why* the result is true. It is always a good idea to look for a proof that illuminates what is going on behind the formula. For this reason, many people prefer the following alternative proof of the proposition:

Proof. Let $f(n) = 0 + 1 + 2 + \cdots + n$. Then, since addition is commutative, we may also write $f(n) = n + (n-1) + (n-2) + \cdots + 1 + 0$. Combining the terms two at a time, we get $2f(n) = [0+n] + [1+(n-1)] + [2+(n-2)] + \cdots + [n+0] = n+n+n+\cdots+n = n(n+1)$. Dividing by two gives the desired result that $f(n) = \frac{n(n+1)}{2}$. \square

Although it appears that induction isn't involved in this second proof, it really is. There are some unproven assertions hidden in the proof, and the proof of these assertions does require proof by induction. Can you find them? You really can't get around induction for proofs of results like this! But you can write the proof in an informal style that enhances understanding. Be prepared to justify your intuitive assertions if asked!

The sum $\sum_{i=0}^n i$ is determined by n alone; in other words, $\sum_{i=0}^n i$ is a *function* of n . For a moment, let us denote this function by f : $f(n) = \sum_{i=0}^n i$. The domain of f is the set of natural numbers, \mathbb{N} . Its outputs can be proven by induction to also be natural numbers, although that is not necessary: the outputs of an inductively defined function can be in any set.

We have defined $f(n)$ for each $n \in \mathbb{N}$ by applying a formula based on an already known function, in this case simply the binary operation of addition, to n and the output of $f(n-1)$: $f(n) = f(n-1) + n$. We refer to this formula as the *recursive rule*. The recursive rule may involve any operations that have already been defined. The inputs for the recursive rule may include, along with the number n itself, any number of previous outputs of the function being recursively defined. Analogous to a proof by induction, you have to start by defining the number of initial outputs needed. In this example, we only needed to define the specific value of $f(0)$ initially. Here is two similar examples:

Example. The definition of the factorial function uses multiplication in its recursive rule rather than addition.

$$n! = \begin{cases} 1, & \text{if } n = 0 \\ n \cdot (n-1)!, & \text{if } n > 0 \end{cases}$$

Informally, $n! = n(n-1)(n-2)\cdots(2)(1)$. If we denote the factorial function by g , then the recursive rule is $g(n) = n \cdot g(n-1)$.

Example. The exponential function for a given base b is defined for natural number inputs as follows:

$$b^n = \begin{cases} 1, & \text{if } n = 0 \\ b \cdot b^{n-1}, & \text{if } n > 0 \end{cases}$$

Any function f whose domain is the natural numbers is equivalent to the infinite sequence $(a_i)_{i \in \mathbb{N}} = (a_i)_{i=0}^{\infty} = (a_0, a_1, a_2, \dots)$ defined by $a_i = f(i)$. For example, the function f defined by $f(i) = i^2$, $i \in \mathbb{N}$, is equivalently defined by the sequence $(i^2)_{i \in \mathbb{N}} = (i^2)_{i=0}^{\infty}$. A nice aspect of sequential notation is that one does not need to designate a letter such as f to name a function defined by a formula. The set of subscripts that indicate the place of each entry in the sequence is called the *index set* (hence the customary use of the letter i , or subsequent letters of the alphabet as needed, to denote an arbitrary element of this set), and its elements are called *indices* (singular: *index*). More generally, any function whose domain is an inductive subset of the natural numbers is equivalent to an infinite sequence. For example, if the domain of a function f is the set $\mathbb{N} \setminus \{0, 1\} = \{2, 3, 4, \dots\}$, then f is equivalently described by the sequence $(a_i)_{i \in \mathbb{N} \setminus \{0, 1\}} = (a_i)_{i=2}^{\infty} = (a_2, a_3, a_4, \dots)$ defined by $a_i = f(i)$. Still more generally, any function at all may be described as a sequence using its domain as the index set. In this way we obtain finite sequences, such as $(i^2)_{i=1}^9$, as well as sequences in which the sequential ordering characteristic of the natural numbers is lost, such as $(x^2)_{x \in \mathbb{R}}$. Sequential notation is just another way to describe a function.

Example. The famous Fibonacci sequence, which begins $(1, 1, 2, 3, 5, 8, 13, \dots)$, is defined by $a_0 = F(0) = 1$, $a_1 = F(1) = 1$, and the recursive rule that $a_n = F(n) = F(n-2) + F(n-1)$ for $n \geq 2$. This is an example of a recursive rule that retrieves output values further back than the preceding value.

It can be proven that any recursive rule and suitable sequence of initial values defines a unique function. The formal statement and proof of this fact are somewhat involved, and we defer them until we have completed a thorough, rigorous discussion of functions in Chapter 6.

5.8 Proving Inequalities by Induction

Proving inequalities is not always as easy as in our first example, that for all natural numbers n , $n \geq 0$. On the one hand, there is more room to maneuver than when proving an equation, since exactitude is not necessary. On the other hand, demonstrable intermediate inequalities sufficient to prove a desired result may not be entirely obvious. Here is another example of an inequality proved by induction:

Proposition. For all natural numbers n , $2^n > n$.

Proof.

Claim (Initial Claim). $2^0 = 1 > 0$.

Claim (Inductive Claim). $2^n > n \Rightarrow 2^{n+1} > n + 1$.

Assume $2^n > n$. Since $2^{n+1} = 2 \cdot 2^n$, by definition, and $2 > 0$, we have $2^{n+1} > 2n$, by Axiom 6(b) and the inductive hypothesis. (By definition, $2 = 1 + 1$, and we previously proved that $1 > 0$; hence $2 = 1 + 1 > 1 + 0 = 1 > 0$.²) Thus, it now suffices to prove that $2n \geq n + 1$. (Observe that \geq is sufficient, since we already have strict inequality in $2^n > 2n$. There is no sense trying to prove something stronger than we have to!) By the distributive law and the defining property of 1, $2n = n + n$. If $n \geq 1$, it is clear (by Axiom 6(a)) that $n + n \geq n + 1$, but it is also possible that $n = 0$. Therefore, we must divide into cases, using a different argument in the case $n = 0$. (We will show independently in the exercises that there are no natural numbers between 0 and 1, so if $n \neq 0$, then $n \geq 1$.)

Case 1: $n = 0$. $2^{n+1} = 2^1 = 2 > 1 = n + 1$. (Observe that the inductive hypothesis is not necessary in this case, since its specificity allows us to deduce the conclusion unconditionally.)

Case 2: $n \geq 1$. As shown above, $2^{n+1} > n + n \geq n + 1$.

□

5.9 Exercises.

1. Prove that if every set in a family contains 0, then the intersection of the sets in this family contains 0.
2. Prove that the intersection of a family of inductive sets is an inductive set.
3. Prove for arbitrary unions and intersections that:

(a) For any family of sets \mathcal{A} , $(\forall A \in \mathcal{A})(A \subseteq \bigcup_{A \in \mathcal{A}} A)$.

(By a *family* of sets we simply mean a set of sets; we choose a different word to make it easier to distinguish verbally the sets that are elements of the family from the family itself. Sometimes other synonyms for *set*, such as *collection*, are used as well.)

²We know from an earlier proposition that $2 \geq 0$. The argument just given rules out the possibility that $2 = 0$. Although this possibility may seem ridiculous, it actually occurs in algebraic systems that don't possess an order satisfying Axiom 6.

(b) For any family of sets \mathcal{A} , $(\forall A \in \mathcal{A})(\bigcap_{A \in \mathcal{A}} A \subseteq A)$.

4. Prove that if S is any inductive set such that $0 \in S$, then $\mathbb{N} \subseteq S$. (Hint: Use the definition of \mathbb{N} as the intersection of all inductive sets containing 0. From a notational point of view, it is necessary to give a name to the collection of all such sets, as in “let \mathcal{I}_0 be the set of all inductive sets containing 0.”)
5. Prove that if m and n are natural numbers, then $n + m$ is a natural number. (Hint: Fix any $m \in \mathbb{N}$ and use induction on n , starting with $n = 0$.)
6. Prove that the product of any two natural numbers is a natural number.
7. Prove that there is no natural number between 0 and 1. (Hint: Show that the set $\{0\} \cup [1, \infty)$ is inductive. The proof naturally divides into two cases.)
8. For every natural number n , prove that there is no natural number between n and $n + 1$.
9. Prove that every natural number other than 0 is greater than or equal to 1.
10. Compute the first ten terms of the Fibonacci sequence.
11. Prove the fundamental properties of exponents:

(a) $b^m b^n = b^{m+n}$.

(b) $(b^m)^n = b^{mn}$.

Hint: For both exercises, fix a natural number m and use induction on n .

12. Let $f : \mathbb{N} \rightarrow \mathbb{R}$ be a fixed function. Define $\sum_{i=0}^n f(i)$ using a recursive definition.
13. Let $f : \mathbb{N} \rightarrow \mathbb{R}$ be a fixed function. Define the general product $\prod_{i=0}^n f(i) = f(0)f(1)f(2) \cdots f(n)$ using a recursive definition.
14. Let r be any fixed real number. Define $f : \mathbb{N} \rightarrow \mathbb{R}$ by $f(i) = r^i$. For this particular function, the sum defined recursively in the manner of the previous problem is called a *geometric sum* and may be written informally as $1 + r + r^2 + \cdots + r^n$. Prove the *geometric sum formula*: $1 + r + r^2 + \cdots + r^n = \frac{1-r^{n+1}}{1-r}$.
15. For all natural numbers n , show that $\sum_{i=1}^n i^2 = \frac{n(n+1)(2n+1)}{6}$.

16. Let a be a fixed real number, and let $f : \mathbb{N} \rightarrow \mathbb{R}$ be a fixed function. . Prove that $a \sum_{i=1}^n f(i) = \sum_{i=1}^n af(i)$.
17. Let $f, g : \mathbb{N} \rightarrow \mathbb{R}$ be fixed functions. Prove that $\sum_{i=1}^n f(i) + g(i) = \sum_{i=1}^n f(i) + \sum_{i=1}^n g(i)$.
18. Combining the results of the previous two problems, conclude that for fixed real numbers a and b and fixed functions $f, g : \mathbb{N} \rightarrow \mathbb{R}$, $\sum_{i=1}^n af(i) + bg(i) = a \sum_{i=1}^n f(i) + b \sum_{i=1}^n g(i)$. (For any real numbers a and b , the function that associates to each ordered pair of real numbers (x, y) the expression $ax + by$ is called a *linear combination* of x and y ; thus, we have proven that “the sum of a linear combination is the linear combination of the sums.”)
19. Use the result of the previous problem to give an alternate (and more intuitive) proof that $1 + r + r^2 + \cdots + r^n = \frac{1 - r^{n+1}}{1 - r}$. (Hint: Consider $(1 - r)(1 + r + r^2 + \cdots + r^n)$.)
20. Apply the trick you used in the previous problem or, alternatively, apply the geometric sum formula, to find the fractional expression for the repeating decimal $\overline{12345}$.
21. Prove, showing all details, that for all natural numbers n , $3^n > n$.
22. Prove in general that if $m \neq 0$ is a natural number, then for all natural numbers n , $m^n > n$.

5.10 Getting the Most out of Induction

5.10.1 Extending to Initial Numbers Other Than Zero

To summarize what we have discussed so far, the method of proof by induction applies to proving that all natural numbers have a given property or, equivalently, that the set of natural numbers satisfying this given property contains all of the natural numbers. The method of proof is to demonstrate that the set of natural numbers satisfying the given property contains 0 and is inductive, since we know that any inductive set of real numbers that contains 0 must contain all the natural numbers.

The following rather obvious result will be useful in the subsequent discussion:

Proposition. *If $n \in \mathbb{N}$ and $n > 0$, then $n - 1 \in \mathbb{N}$.*

Proof. We can rephrase the proposition as follows (using one of our previous results):
 $\forall n \in \mathbb{N}, n = 0$ or $n - 1 \in \mathbb{N}$.

Obviously, the proposition is true for 0.

Suppose it is true for n . We must prove it is true for $n + 1$. Since it is true for n , either $n = 0$ or $n - 1 \in \mathbb{N}$.

Case 1: $n = 0$. Then $(n + 1) - 1 = 0 \in \mathbb{N}$.

Case 2: $n - 1 \in \mathbb{N}$. Then $(n + 1) - 1 = (n - 1) + 1 \in \mathbb{N}$, since $n - 1 \in \mathbb{N}$ and \mathbb{N} is inductive.

□

Some properties are only true for natural numbers beyond a certain point. For example, $2^n > n^2$ for every $n \geq 5$, but it is not true for 2, 3, or 4. A simple extension of proof by induction allows us to prove propositions such as this. We simply begin our induction at the appropriate starting point, instead of at 0. That is, to prove that $(\forall n \in \mathbb{N} : n \geq N)p(n)$, we prove that $\{n \in \mathbb{N} : n \geq N \text{ and } p(n)\}$ contains N and is inductive. As an example, consider the proposition just mentioned, in which $N = 5$:

Proposition. For every $n \geq 5$, $2^n > n^2$.

Proof. We first check that $2^5 = 32 > 25 = 5^2$.

Now suppose $2^n > n^2$. We must prove that $2^{n+1} > (n+1)^2$. Now, $2^{n+1} = 2 \cdot 2^n > 2n^2$ by inductive hypothesis. (Here we also use Axiom 6(b), together with the fact that $2 > 0$, which follows as a particular case from the previously proven proposition that every natural number is greater than 0. That 2 is, in fact, a natural number follows from its definition and the inductive property of the natural number system: $2 = 1 + 1$. All of our familiar numerals for natural numbers are defined inductively using the rules of the base ten place value system: $3 = 2 + 1$, $4 = 3 + 1$, etc.) Since our goal is to prove $2^{n+1} > (n+1)^2$, it now suffices to prove $2n^2 > n + 1^2$. Since $(n+1)^2 = n^2 + 2n + 1$, it suffices to prove that $n^2 > 2n + 1$. Given that $n \geq 5 > 1$, we have $n^2 \geq 5n = 2n + 3n$, and $3n > 1$. □

This method works because $\{n \in \mathbb{N} : n \geq N\}$ is the smallest inductive set containing N (as we will prove); therefore, $\{n \in \mathbb{N} : n \geq N\}$ must be a subset of any inductive set that contains N . More generally, if N is any integer, $\{n \in \mathbb{Z} : n \geq N\}$ is the smallest inductive set containing N , so an inductive proof can begin with negative integers, too. (Recall that \mathbb{Z} is the set consisting of all natural numbers and their additive inverses: $\mathbb{Z} = \mathbb{N} \cup \{-n : n \in \mathbb{N}\}$.) Alternatively, we could view the proposition above as being

conditional: $(\forall n \in \mathbb{N})(n \geq 5 \Rightarrow 2^n > n^2)$; hence, ordinary induction initiated with 0 is valid (and the proposition holds for $n < 5$ by virtue of they condition being false). A disadvantage of this justification, however, is that it does not extend to initiation with negative integers.

We will work in small manageable steps toward the proof that, for any integer N , $\{n \in \mathbb{Z} : n \geq N\}$ is the smallest inductive set containing N . First we assert that the set of integers is inductive:

Lemma. *If $n \in \mathbb{Z}$, then $n + 1 \in \mathbb{Z}$.*

Proof. The lemma is true for all $n \in \mathbb{N}$ by definition of \mathbb{N} , so we may restrict our attention to integers of the form $-n$, where $n \in \mathbb{N}$ and $n > 0$. We have that $-n + 1 = -(n - 1)$, and by a previous proposition, if $n > 0$, then $n - 1 \in \mathbb{N}$; therefore, $-n + 1 = -(n - 1) \in \mathbb{Z}$. \square

Next we assert that the set $\{n \in \mathbb{Z} : n \geq N\}$ is inductive:

Lemma. *$\{n \in \mathbb{Z} : n \geq N\}$ is inductive.*

Proof. Since \mathbb{Z} is inductive, by the previous lemma, we need only prove that if $n \geq N$, then $n + 1 \geq N$. The argument, easily constructed using the fact that $1 > 0$, is left to the reader. \square

Next we make the following observation about differences of natural numbers:

Lemma. *For all natural numbers m and n , if $m \leq n$, then $n - m \in \mathbb{N}$.*

Proof. Let $n \in \mathbb{N}$. We prove the lemma by induction on m .

Claim. *If $m \geq 0$, then $n - 0 \in \mathbb{N}$.*

This is trivial, since $n \in \mathbb{N}$ by hypothesis.

Claim. *If the result holds for m , then it holds for $m + 1$.*

Assume the result holds for m and that $n \geq m + 1$. Since $m + 1 > m$, clearly $n > m$. Since the result holds for m , $n - m \in \mathbb{N}$. Furthermore, $n - m > 0$. It follows from a previous proposition that $n - (m + 1) = (n - m) - 1 \in \mathbb{N}$.

We have proven that the set of natural numbers for which the lemma holds is inductive and contains 0; therefore, the lemma holds for all natural numbers. \square

Remark. Alternatively, we could have proved the proposition by induction on n , which we encourage the reader to try for practice. The initial claim is trivial, since no natural number is less than 0. The inductive claim is that if for all $m \leq n$, $m - n \in \mathbb{N}$, then for all $m \leq n + 1$, $n + 1 - m \in \mathbb{N}$. To prove it, divide into the cases that $m \leq n$ and $m = n + 1$.

Lastly before proving the proposition, we use the lemma above to prove the following more general observation:

Lemma. *If $m \in \mathbb{Z}$, $n \in \mathbb{Z}$, and $m \leq n$, then $n - m \in \mathbb{N}$.*

Proof. By the previous lemma we may restrict our attention to the case that $m = -k$, where $k \in \mathbb{N}$.

Case 1: $n \in \mathbb{N}$. Then $n - m = n + k \in \mathbb{N}$ by the result of Exercise 5 of the previous section.

Case 2: $n = -l$, where $l \in \mathbb{N}$. Since $m \leq n$ by hypothesis, $-k \leq -l$, from which it follows that $k \geq l$. Thus $n - m = -l + k = k - l \in \mathbb{N}$ by the result of the previous lemma.

□

Proposition. *Let N be an integer, and let \mathcal{I}_N be the set of all inductive sets containing N . Then $\{n \in \mathbb{Z} : n \geq N\} = \bigcap_{A \in \mathcal{I}_N} A$.*

Proof. We have proven above that the set $\{n \in \mathbb{Z} : n \geq N\}$ is inductive; hence, $\{n \in \mathbb{Z} : n \geq N\} \in \mathcal{I}_N$. Therefore, $\bigcap_{A \in \mathcal{I}_N} A \subseteq \{n \in \mathbb{Z} : n \geq N\}$ by the result of Exercise 3 of the previous section.

Conversely, let $n \in \{n \in \mathbb{Z} : n \geq N\}$.

By the previous lemma, $n = N + m$ for some $m \in \mathbb{N}$. We complete the proof by showing that $N + m \in \bigcap_{A \in \mathcal{I}_N} A$ for all natural numbers m .

Claim. $N + 0 \in \bigcap_{A \in \mathcal{I}_N} A$.

By definition of \mathcal{I}_N , N is an element of every set $A \in \mathcal{I}_N$; thus, by definition of intersection, $N + 0 \in \bigcap_{A \in \mathcal{I}_N} A$.

Claim. *If $N + m \in \bigcap_{A \in \mathcal{I}_N} A$, then $N + m + 1 \in \bigcap_{A \in \mathcal{I}_N} A$.*

By hypothesis and the definition of intersection, $N + m$ is an element of every set $A \in \mathcal{I}_N$. Since every set $A \in \mathcal{I}_N$ is inductive, by definition of the collection \mathcal{I}_N , $N + m + 1$ is an element of every set $A \in \mathcal{I}_N$; hence, by definition of intersection, $N + m + 1 \in \bigcap_{A \in \mathcal{I}_N} A$. □

5.10.2 Complete Induction: Two Approaches

Modifying the proposition to enhance the inductive step

We have seen how to extend proof by induction by using an appropriately modified initial step. It is also possible to greatly increase the power of proof by induction by refining the inductive step. These two enhancements may be combined, but for simplicity let us restrict our attention to proving a proposition $p(n)$ for all natural numbers n .

To see how we might improve our situation with respect to the inductive step, consider a property $q(n)$ that is in general *stronger* than $p(n)$, meaning that $q(n) \Rightarrow p(n)$ for all $n \in \mathbb{N}$. We could obtain the desired result that $(\forall n \in \mathbb{N})p(n)$ by proving the stronger result that $(\forall n \in \mathbb{N})q(n)$ instead. At first glance, this might seem counter-intuitive. Wouldn't it be *harder* to prove something stronger? Not necessarily! Remember that in the inductive step, we would not be proving $q(n)$ directly, but rather the *conditional* statement $q(n) \Rightarrow q(n+1)$. Thus, a stronger property $q(n)$ gives us a stronger inductive hypothesis to work with! Of course, $q(n+1)$ is a stronger inductive conclusion, but there is a clever way to choose $q(n)$ so that everything additional in the conclusion follows immediately from the stronger hypothesis. In addition, $q(0)$ will be equivalent to $p(0)$, so no more work required for the initial step.

Here's how! Consider the statement $q(n) = (\forall i \in \mathbb{N} : i \leq n)p(i)$ (which is clearly stronger than $p(n)$). Then $q(n) \Rightarrow q(n+1)$ translates as $(\forall i \in \mathbb{N} : i \leq n)p(i) \Rightarrow (\forall i \in \mathbb{N} : i \leq n+1)p(i)$. In other words, our new inductive hypothesis is that the property given by $p(i)$ is true for all natural numbers that are less than or equal to n , and the new inductive conclusion is that it is true for all natural numbers less than or equal to $n+1$. Now, since the inductive hypothesis already tells us directly that the property $p(i)$ is true for all natural numbers less than or equal to n (and there are no natural numbers between n and $n+1$), the only new assertion in the inductive conclusion is that it is true for $n+1$. This is just the same conclusion we had before! The conclusion is the same as before, but the hypothesis is much stronger, since you can use the truth of the property for *all* the numbers up to and including n , not just n itself! Furthermore, since there are no natural numbers less than 0, the initial step is also the same.

The technique of using the individual assertion $p(n)$ as the inductive hypothesis is sometimes called *simple* induction, in contrast to using the more general, and therefore stronger, assertion that $(\forall i \in \mathbb{N} : i \leq n)p(i)$, which is sometimes called *complete* induction. (You are using the complete power that induction gives you.) Informally, you can think of simple induction (starting at 0 for simplicity) as: "It is true for 0. If it is true for 0, then it must be true for 1. If it is true for 1, then it must be true for 2. If it is true for 2, then it must be true for 3. . . . Therefore, it is true for all natural numbers."

You can think of complete induction as: “It is true for 0. If it is true for 0, then it must be true for 1. If it is true for 0 and 1, then it must be true for 2. If it is true for 0, 1, and 2, then it must be true for 3. . . . Therefore, it is true for all natural numbers.”

Often, simple induction is sufficient, but sometimes complete induction is absolutely necessary. A classic example is the following fundamental result of number theory. To understand its statement, we must recall the following definitions:

Definition. The natural number m *divides* n if there is some natural number k such that $km = n$. We also say that m is a *factor* of n .

Notation. Denote the fact that m divides n by $m|n$.

Clearly, every natural number has 1 and itself as factors. The numbers that have no other factors are special.

Definition. A natural number p is *prime* if it is greater than 1 and its only factors are 1 and p . A natural number that is not prime is called *composite*.

Note that a factor of a number cannot be larger than the number. (Why?) Therefore, to check that a number is prime, it is only necessary to check that the smaller natural numbers other than 1 do not divide it. A quick check demonstrates, for example, that 2 and 3 are prime numbers. (In fact, you can be more efficient than that; it is only necessary to check numbers less than or equal to the square root of the number. Why?)

Theorem. *Any natural number greater than 1 can be factored as the product of prime factors.*

Example. $24 = 2 \cdot 2 \cdot 2 \cdot 3$.

Proof. The initial step is to show that 2 factors into primes, which is trivial, since 2 itself is prime.

For the inductive step, let n be any natural number greater than 1. Assume that any number less than n factors into primes.

Case 1: n is prime. Then we are done.

Case 2: n is composite. Then, by definition, $n = lm$ for some natural numbers l and m , neither of which is 1. By inductive hypothesis, l and m factor as the product of prime factors, say $l = p_1 p_2 \cdots p_j$ and $m = p_{j+1} p_{j+2} \cdots p_k$. (Note that the prime factors p_i in these expressions need not be distinct.) Then $n = p_1 p_2 \cdots p_j p_{j+1} p_{j+2} \cdots p_k$ is the product of prime factors.

□

Remark. Simple induction will not work for this theorem, since neither factor of n will be $n - 1$. The simple inductive hypothesis that $n - 1$ factors into primes is not merely insufficient, but useless.

Remark. It is also true that, up to the order of the factors, the prime factorization of a number is unique. The uniqueness of the factorization is also proved by complete induction, but the proof is considerably more complicated.

Well-ordering

An approach that is both more elegant and more general in its applications is to introduce the concept of *well-ordering*. First we define the smallest element of a totally ordered set in the obvious way:

Definition. Let S be an ordered set. An element $s_0 \in S$ is the *smallest* element of S if for every $s \in S$, $s_0 \leq s$.

It is easy to prove that a set can have at most one smallest element, justifying the anticipatory use of the word “the” in the definition. We leave this proof as an exercise.

A set need not have a smallest element. The following are examples of sets that do not have a smallest element: \emptyset , \mathbb{Z} , the interval $(0, 1)$, the interval $(0, 1]$. The empty set is a special case: it fails to have a smallest element simply because it has no elements at all.

Definition. An ordered set T is *well-ordered* if every non-empty subset $S \subseteq T$ has a smallest element.

Remark. Since the empty set is a subset of every set, we must exclude this special case in the definition, or it would be vacuous: no well-ordered set would exist.

Remark. We need not assume in the definition that the ordering of T is total, but in fact it must be. Consider any two distinct elements t and t' ; since the set $\{t, t'\}$ has a smallest element, one must be less than the other.

Not every totally ordered set is well-ordered, by any means. The sets \mathbb{Z} , $(0, 1)$, and $(0, 1]$ are not well-ordered because they do not have a smallest element. The intervals $[0, 1)$, $[0, 1]$, and $[0, \infty)$ are not well-ordered, even though each of them does have a smallest element, because many of their non-empty subsets, such as $(0, 1)$, do not. The set $\{0\} \cup \{\frac{1}{n} : n \in \mathbb{N}\}$ is not well-ordered because every infinite subset that does not contain 0 fails to have a smallest element.

Well-ordered sets are very special, and of course we have a reason for defining them. The following theorem leads to a powerful method of proving universal propositions about their elements:

Theorem (Principle of Complete Induction). *Let T be a well-ordered set. Let $p(t)$ be a proposition about the elements of T such that $(\forall t \in T : t < s)p(t) \Rightarrow p(s)$. Then $(\forall t \in T)p(t)$.*

Proof. Consider $S = \{t \in T : \neg p(t)\}$. We will show $S = \emptyset$; it follows that $p(t)$ holds for all $t \in T$. Suppose to the contrary that S is not empty. Then since T is well-ordered, S has a smallest element, s_0 . Thus for all $t \in T$ such that $t < s_0$, $t \notin S$, and hence $p(t)$ holds. But then, by our assumption with respect to $p(t)$, $p(s_0)$ holds, contradicting the assertion that $s_0 \in S$. \square

Many well-ordered sets do exist, so the Principle of Complete Induction has wide application. As we know, the method of complete induction is valid for the set of natural numbers, so it should come as no surprise that \mathbb{N} , with its standard ordering (as a subset of \mathbb{R}), is well-ordered. In fact, a set T is well-ordered if and only if the method of complete induction is valid for T :

Theorem. *A set T is well-ordered if and only if for any proposition $p(t)$ about the elements of t such that $(\forall t \in T : t < s)p(t) \Rightarrow p(s)$, $(\forall t \in T)p(t)$.*

Proof. The “only if” part is the Principle of Complete Induction, proven above. For the “if” part we will prove that any subset of T that fails to have a smallest element is empty. Let $S \subseteq T$ fail to have a smallest element, and consider the proposition that $t \notin S$. If $t \notin S$ for all $t < s$, then surely $s \notin S$, since if s were an element of S , it would be the smallest! Thus, by hypothesis, $t \notin S$ for all $t \in T$; hence S is empty. \square

Corollary. *The set of natural numbers, with its standard ordering, is well-ordered.*

Remark. No algebraic operations are used in these definitions, theorems, or proofs. A well-ordered set need not be a subset of the set of real numbers, nor need it have any operations defined on it; the method of complete induction will apply to it nonetheless. Conversely, a set for which the method of complete induction is valid need not be a subset of the set of real numbers or have any operations defined on it; it will be well-ordered nonetheless. (Of course, propositions about the elements of the set can only use operations and relations that are defined on it.)

As outlined in the exercises, one can prove directly by simple induction that \mathbb{N} is well-ordered, obtaining as a consequence the validity of complete induction.

5.11 Exercises

1. Prove that the sum of any two integers is an integer.
2. Prove that the product of any two integers is an integer.
3. Prove that for all natural numbers $n \geq 4$, $2^n < n!$.
4. Prove that for all natural numbers n , $3^n > n^2$.
5. Prove that for all natural numbers $n \geq 10$, $2^n > n^3$.

Remark. It is true in general that, for any natural number $m \neq 0$ and for any fixed exponent $k \in \mathbb{N}$ and for some sufficiently large $N \in \mathbb{N}$, $n \geq N \Rightarrow m^n > n^k$. However, the algebra involved in the method of proof we have used becomes unwieldy for larger values of k . The methods of calculus are more effective for proving these results.

6. Prove that, for any natural number $n \geq 12$, there are natural numbers k and l such that $n = 3k + 7l$. (Hint: You will need complete induction. Check 13 and 14 separately; for the inductive step, consider $n - 3$. Why must you check 13 and 14 separately?)
7. Prove that a set can have at most one smallest element.
8. This exercise outlines a direct proof using simple induction that \mathbb{N} is a well-ordered set. Complete the proofs of each of the following claims, supplying all details and justifications.
 - (a) For any natural number n , the set of natural numbers less than or equal to n (informally, $\{0, 1, 2, \dots, n\}$) is well-ordered.

Hints: Use proof by induction. The initial claim is that $\{0\}$ is well-ordered, since its only non-empty subset, namely $\{0\}$ itself, has a smallest element, namely 0, its sole element. For the inductive claim, assume that $\{0, 1, 2, \dots, n-1\}$ is well-ordered. To prove that $\{0, 1, 2, \dots, n\}$ is well-ordered, let S be a non-empty subset of $\{0, 1, 2, \dots, n\}$. Consider $S \cap \{0, 1, 2, \dots, n-1\}$, and divide into two cases:

Case 1: $S \cap \{0, 1, 2, \dots, n-1\} = \emptyset$, in which case $S = \{n\}$ (recall that we proved there are no natural numbers between $n-1$ and n), or

Case 2: $S \cap \{0, 1, 2, \dots, n-1\} \neq \emptyset$, in which case $S \cap \{0, 1, 2, \dots, n-1\}$ is a non-empty subset of $\{0, 1, 2, \dots, n-1\}$ and hence has a smallest element by the inductive hypothesis.

(b) \mathbb{N} is well-ordered (the desired full result).

Hints: Let S be any non-empty subset of \mathbb{N} . Since $S \neq \emptyset$, there is a natural number $n \in S$; therefore, $S \cap \{0, 1, 2, \dots, n\}$ is a non-empty subset of $\{0, 1, 2, \dots, n\}$. Now use the result of part (a), noting that any element of S that is not a member of $S \cap \{0, 1, 2, \dots, n\}$ is greater than n .

9. (a) Prove that for any natural number n and any natural number m other than 0, there is a number l such that $lm > n$. (Hint: Let $l = n + 1$. This is a sledgehammer approach, $n + 1$ generally being a much larger value of l than is necessary, but it works!)
- (b) Prove the existence of an algorithm for division with remainder in the natural number system: for any natural numbers m and n such that $0 < m \leq n$, there exist natural numbers q (the *quotient*) and r (the *remainder*) such that $n = qm + r$ and $r < m$. (Note that $q \neq 0$, since $r < m$.)

Hints: Use the fact that \mathbb{N} is well-ordered: Consider the set $S = \{l \in \mathbb{N} : lm > n\}$. By the result of part (a), $S \neq \emptyset$; therefore, it has a smallest element. Let k be the smallest element of S . Clearly $0 \notin S$; hence, $k > 0$. (Furthermore, since $m \leq n$, $1 \notin S$; hence, $k > 1$, although this fact is not needed for the proof). Thus $k - 1 \in \mathbb{N}$. Let $q = k - 1$.

Remark. This theorem does not actually provide a practical algorithm: considering all multiples of the divisor until one found the very largest that is less than the dividend could be quite time consuming! The familiar long division algorithm simplifies the procedure by taking advantage of place-value numeration.

10. This exercise outlines a proof that any two natural numbers, m and n , have a unique greatest common divisor, $GCD(m, n)$, which is not only the largest but also a multiple of any other common divisor. Furthermore, as a corollary of the proof³, we obtain the invaluable result that $GCD(m, n)$ is a linear combination of m and n with integer coefficients (that is, there are integers k and l such that $GCD(m, n) = km + ln$).

³The technical term for a corollary of the proof, as opposed to a corollary of the theorem itself, is a *skolem*.

- (a) Let $S = \{km + ln : k, l \in \mathbb{Z} \text{ and } km + ln > 0\}$. Prove that $S \neq \emptyset$.
- (b) Therefore, since \mathbb{N} is well-ordered, S has a smallest element. Let d be the smallest element of S . Prove that any common divisor of m and n divides d .
- (c) Prove that $d \leq m$ and $d \leq n$.
- (d) Prove that d itself divides both m and n . (Hints: By Exercise 9b above, $m = qd + r$, where $0 \leq r < d$. The remainder r is a linear combination of m and n , since $r = m - qd$; however, since $r < d$, and d is by definition the smallest element of S , $r \notin S$. Thus the only possibility is that $r = 0$. The proof that $d|n$ is similar.)
- (e) Prove that if d' is a natural number such that $d'|m$, $d'|n$, and any common divisor of m and n divides d' , then $d' = d$. Thus we are justified in defining $GCD(m, n)$ to be the unique natural number with these properties. It also follows from the definition of d that there are integers k and l such that $d = km + ln$.

The following exercises lead you through the proof that prime factorizations are unique. Prove each of the propositions.

11. If p is prime and $p|mn$, then $p|m$ or $p|n$. (Hint: Suppose $p \nmid m$. Then since p is prime, $GCD(p, m) = 1$. By the previous result, there exist integers k and l such that $1 = kp + lm$. Use this and the fact that $n = 1n$ to show $p|n$.)
12. More generally, if p is prime and $p|m_1m_2m_3 \cdots m_n$, then $p|m_i$ for some i . (Hint: The case $n = 2$ is the result of the previous exercise. Use induction on n .)
13. If $m = p_1p_2p_3 \cdots p_n$ is a factorization of m into primes and $p|m$, then $p = p_i$ for some i .
14. The factorization of any natural number into primes is unique. (Hint: Use complete induction, starting with the number 2.)

The result that every natural number greater than 1 has a unique factorization into prime factors is called the Fundamental Theorem of Arithmetic on account of the depth and breadth of its ramifications. In particular, since rational numbers are ratios of integers, the Fundamental Theorem of Arithmetic has profound consequences for the rational number system. For one thing, we may assume that the numerator and denominator of the fraction representing a rational number have no common factors: just factor into primes and cancel. This simple fact provides the basis for proving that many real numbers are irrational. Prove each of the following propositions.

15. $\sqrt{2}$ is irrational. In other words, there exist no natural numbers m and n such that $\left(\frac{m}{n}\right)^2 = 2$. (As of yet, since we have not examined the continuity of the real number system, we do not know there is any real number with this property! But if there is one (and indeed there is), it cannot be rational.)
16. $\sqrt{3}$ is irrational.
17. More generally, if n is not a perfect square, then \sqrt{n} is irrational. (We say a natural number n is a *perfect square* if there is a natural number m such that $n = m^2$. Hint: Prove the contrapositive. Using the Fundamental Theorem of Arithmetic, first prove the lemma that $q^2|p^2 \Leftrightarrow q|p$.)
18. $\sqrt{2}$ and $\sqrt{3}$ are incommensurable.

Chapter 6

Relations

Mathematics is about relationships. Proofs describe the logical relationship among mathematical assumptions and propositions. Sets describe the association of mathematical objects that play a particular role, such as numbers or points. There are usually additional relationships among the elements of a set or between elements of different sets. Some numbers are larger than other numbers; some points are between pairs of other points. For any pair of points in a geometric object, a certain number, and only that number, is the distance between them. For every pair of numbers, a certain number is their product and a certain number (generally different, but not always) is their sum. For any number, a certain number is its square, and two specific numbers are its square roots. (The latter holds even for negative numbers if we are working in the complex number system). Some pairs of natural numbers have the same parity, either both even or both odd, and other pairs of numbers do not.

Relationships of this type are described formally by set-theoretic constructions called *relations*. Relations lie at the very heart - and brain - of mathematics. Without the concept of a relation, none of the mathematics with which you are familiar would formally exist, not even counting, since counting requires proceeding from each number to its successor; the association between each number and its successor is a relation (more specifically, a function)!

We have already considered the relation of order on the real numbers, as well as those that associate pairs of numbers to their sums and products, and we were able to write these relations down, describe their properties precisely, and prove propositions that involve them. But what exactly *are* they as mathematical objects?

6.1 Ordered pairs

To describe a relation, we must describe which elements of which sets are related, but merely grouping them into sets of two is not enough. We must also be able to describe *how* they are related. For example, 2 is the successor of 1, but 1 is not the successor of 2. Relations are directed. Describing a directed pair requires that we can distinguish one element as coming first and the other as coming second. Sets do not do this. Two sets are equal, by definition, if and only if they contain the same elements. $\{1, 2\} = \{2, 1\}$. Thus we need a new type of object called an *ordered pair*. You are undoubtedly familiar with ordered pairs, but you are likely non familiar with how to define them using sets. The following exercises guide you to the correct definition. (One could certainly think of ordered pairs as primary, undefined objects like sets and simply postulate that the ordered pair (a, b) is equal to the ordered pair (c, d) if and only if $a = c$ and $b = d$; however, defining ordered pairs in terms of sets is a useful exercise that provides insight into the nature of sets themselves.)

6.2 Exercises

- Which of the following pairs of sets has the property that the two sets are equal if and only if $a = c$ and $b = d$? (There is only one! This is trickier than it looks!)
 - $\{a, b\}, \{c, d\}$
 - $\{a, \{b\}\}, \{c, \{d\}\}$
 - $\{\{a\}, b\}, \{c, \{d\}\}$
 - $\{\{a\}, \{a, b\}\}, \{\{c\}, \{c, d\}\}$
 - $\{\{a\}, \{b\}\}, \{\{c\}, \{d\}\}$
 - $\{a, b, \{b\}\}, \{c, d, \{d\}\}$
 - $\{a, b, \{a, b\}\}, \{c, d, \{c, d\}\}$
- It is obvious for each pair of sets in exercise (1) that if $a = c$ and $b = d$, then the sets are equal. As noted, the converse is true for only one pair. Prove the converse for this pair, and give a counterexample to the converse for each of the others.
- Fill in the blank: **Definition.** Given elements a and b of sets A and B , respectively, the *ordered pair* (a, b) is the set _____.

6.3 The Cartesian Product Operation

We define the *Cartesian product* of two sets as follows:

Definition. Given two sets A and B , the *Cartesian product* $A \times B$ is the set $\{(a, b) : a \in A \wedge b \in B\}$.

Note that the order in which the two sets appear matters. It is ambiguous to say “the Cartesian product of sets A and B ,” since the logical connective “and” is symmetric. (Nonetheless, people do say that, with the understanding that A comes first.)

Remark. The Cartesian product operation is named after René Descartes, who introduced the idea of labeling the points of a Euclidean plane with pairs of real coordinates. (The pairs must be ordered, since going 5 to the right and 2 up, say, is different from going 5 up and 2 to the right.)

Remark. Although we are not yet equipped to formally prove this, you should be able to see that if A is a finite set with m elements and B is a finite set with n elements, then $A \times B$ is finite and has mn elements.

6.4 Exercises

1. Prove that $A \times B = \{(a, b) : a \in A \text{ and } b \in B\}$ is a set in accordance with the axioms of set theory. (Hints: By the Axiom of Restricted Comprehension, it suffices to show that $A \times B$ is a subset of a known set. Noting that both elements of each ordered pair are subsets of $A \cup B$, use the Axiom of the Power Set twice.)
2. Prove $(A = \emptyset \text{ or } B = \emptyset) \Leftrightarrow A \times B = \emptyset$. (Hint: Prove the contrapositive biconditional. Suppose there is an element in $A \times B$, and show how it gives rise to the existence of an element in A and an element of B . Conversely, show how the existence of an element of A and an element of B gives rise to an element of $A \times B$.)
3. Prove that if $A \neq \emptyset$ and $B \neq \emptyset$ and $A \neq B$, then $A \times B \neq B \times A$.
4. Prove that if $A \subseteq C$ and $B \subseteq D$, then $A \times B \subseteq C \times D$.
5. Is the converse to the previous proposition true? If so, prove it. If not, give a counterexample. In addition, if the converse is false, find an additional condition on A and B under which it is true, and prove it.

6. Prove that if $A \subseteq C$, $B \subseteq D$, and either $A \neq C \wedge D \neq \emptyset$, or $B \neq D \wedge C \neq \emptyset$, then $A \times B \subsetneq C \times D$. (The symbol \subsetneq means “is a proper subset of,” as should be evident. Sometimes the symbol \subset is used to mean this, but usage varies among authors; \subsetneq leaves no doubt.)
7. Is the converse to the previous proposition true? If so, prove it. If not, give a counterexample. In addition, if the converse is false, find an additional condition on A and B under which it is true, and prove it.

For each of the following propositions, determine if the equation is true. If it is, prove it. If it is not, determine whether or not inclusion holds in one direction. Prove any inclusion that holds, and give a counterexample to any that doesn't.

Finally, for *extra credit*, for any equation that does not hold, if you can find additional conditions under which it does hold, do so and prove your result!

(Note that “ \setminus ” means set difference; it is synonymous with “ $-$.”)

8. $(A \times B) \cup (C \times D) = (A \cup C) \times (B \cup D)$. (Hint: Draw a picture in the coordinate plane!)
9. $(A \times B) \cap (C \times D) = (A \cap C) \times (B \cap D)$.
10. $A \times (B \setminus C) = (A \times B) \setminus (A \times C)$.
11. $(A \setminus B) \times (C \setminus D) = (A \times C) \setminus (B \times D)$.

6.5 The Definition of a Relation

Formally, a *relation* is simply a subset of the Cartesian product of two sets.

Definition. A *relation* of set A to set B is a subset of $A \times B$.

(Note the asymmetry of the language “of ... to ...,” to indicate that A is the first factor of the Cartesian product and B the second.)

This definition is enormously powerful in its simplicity and the scope of its application. The function that associates each real number to its square is just the set of pairs $\{(x, y) \in \mathbb{R} \times \mathbb{R} : x^2 = y\}$. To describe (or postulate) an order relation on a set A , simply describe (or postulate) a subset of $A \times A$ satisfying the appropriate properties; $x < y$ if and only if (x, y) is in this subset. If S is a set of points in which distances can be measured, these distances are completely described by a subset of $(S \times S) \times \mathbb{R}$ satisfying

appropriate properties, associating a unique positive real number to each pair of elements of S . In short, to specify a relation, we just specify the set of pairs of elements that are related.

Many relations relate elements of the same set. Since these are so common, we have special terminology for a relation of a set to itself.

Definition. A *relation of on set* A is a relation of A to itself, that is, a subset of $A \times A$.

(Note that although we may now use symmetric language to describe the factors of the Cartesian product, as they are the same, the relation itself need not be symmetric. For example, an order relation on a set is not symmetric.)

6.6 Operations on Relations: Inverse and Composition

Given a relation R of A to B , it is natural to consider the relation of B to A given by reversing the order of the pairs. This latter relation is called the *inverse* of R and denoted by R^{-1} .

Definition. Given a relation R of A to B , the *inverse* relation, denoted by R^{-1} , is the relation of B to A defined by $R^{-1} = \{(b, a) \in B \times A : (a, b) \in R\}$.

For example, the inverse to the $<$ relation on \mathbb{R} is the $>$ relation on \mathbb{R} . As another example, the inverse of the relation of points lying on lines is the relation of lines passing through points. (Note that inverse relations have nothing to do with taking reciprocals.)

There is also a natural method of combining two relations. If R relates set A to set B , and S relates set B to set C , the *composition* $S \circ R$ is the relation of A to C that pairs those elements associated with a common element of B by R and S , respectively. (You might wonder why S is written on the left, when it would seem natural to write it on the right. The reason is to make this notation compatible with the standard functional notation, described below. The operation of composition is typically used with functions.)

Definition. Given a relation R of A to B and a relation S of B to C , the *composition* $S \circ R$ is the relation of A to C defined by $S \circ R = \{(a, c) \in A \times C : (\exists b \in B) [(a, b) \in R \wedge (b, c) \in S]\}$.

Some special types of relations are so pervasive in mathematics that they merit separate, detailed discussion. These are functions, order relations, and equivalence relations,

the topics of the following sections. I suspect you are at least informally familiar with these notions. However, keep in mind that not all important relations fit into one of these categories. The relationship of betweenness for points on a line is not a function, order relation, or equivalence relation. (Once coordinates are chosen on a line, betweenness can be described in terms of the order of these coordinates, but is not itself an order relation.) Neither is the relation of *incidence* between points and lines. (“Incidence” is the technical term. We more commonly say that a point “lies on” a line, or that a line “passes through” a point.) The association of a number to its square roots is the inverse relation to a function, but is not itself a function (or order relation or equivalence relation).

6.7 Exercises

- Let R be the relation on \mathbb{R} given by $\{(x, y) : x < y\}$. Then R^{-1} is the relation on \mathbb{R} given by $\{(x, y) : y < x\}$. Describe:
 - $R^{-1} \circ R$
 - $R \circ R^{-1}$
 - $R \circ R$
- Let S be the relation on \mathbb{R} given by $\{(x, y) : x^2 = y\}$. Describe:
 - S^{-1}
 - $S^{-1} \circ S$
 - $S \circ S^{-1}$
 - $S \circ S$.
- Let T be the relation on \mathbb{R} given by $\{(x, y) : x = \pm y\}$. Describe:
 - T^{-1}
 - $T^{-1} \circ T$
 - $T \circ T^{-1}$
 - $T \circ T$.
- Let U be the relation on \mathbb{R} given by $\{(x, y) : x = \pm y \wedge |x| \leq 1\}$. (We define $|x| = x$, if $x \geq 0$, and $|x| = -x$, if $x < 0$.) Describe:

- (a) U^{-1}
 - (b) $U^{-1} \circ U$
 - (c) $U \circ U^{-1}$
 - (d) $U \circ U$.
5. With R , S , T , and U as defined above, describe:
- (a) $R \circ S$
 - (b) $S \circ R$
 - (c) $R \circ T$
 - (d) $T \circ R$
 - (e) $R \circ U$
 - (f) $U \circ R$
 - (g) $S \circ T$
 - (h) $T \circ S$
 - (i) $S \circ U$
 - (j) $U \circ S$
 - (k) $T \circ U$
 - (l) $U \circ T$
6. One may graph any relation on \mathbb{R} in the coordinate plane. Graph each of the relations in the preceding exercises.
7. Let $R \subseteq A \times B$ be a relation. Prove that $(R^{-1})^{-1} = R$. Thus, inverse relations come in pairs, just like additive and multiplicative inverses of numbers.
8. Let R be a relation of A to B and S be a relation of B to C . Prove that $(S \circ R)^{-1} = R^{-1} \circ S^{-1}$.

6.8 Functions

6.8.1 The Definition of a Function

A *function* is a relation that associates a unique second coordinate, which we think of as the “output,” to each first coordinate, thought of as the “input.” Sometimes the

term “independent variable” is used for the input, and the term “dependent variable” is used for the output, since it depends on the input. The latter terminology is somewhat misleading, however, because the second coordinate of any relation depends on the first, it just need not be uniquely determined by it. For example, if I ask which real numbers are greater than, say, 1, my answer depends on the number 1. The set of numbers satisfying this condition is infinite, but not arbitrary, and the set of numbers greater than, say, 2 is different from the set of those greater than 1. The key defining property of a function is that the first coordinate determines the second coordinate *uniquely*. Another way to look at this is that each element of the first factor appears in only one pair in the relation.

Definition. A *function from A to B* is a relation of A to B satisfying the additional property that each element of A occurs as the first coordinate of one and only one element of the relation. If $a \in A$, then the unique element $b \in B$ for which (a, b) is an element of the function is called the *image* of a . The first factor, A , is called the *domain* of the function. The second factor, B , is called the *range* of the function. The subset of the range whose elements actually occur as second coordinates of ordered pairs in the function (that is, images of elements in the domain) is called the *image* of the function.

Notation. If f is a function from A to B , we write $f : A \rightarrow B$. If $(a, b) \in f$, we write $b = f(a)$.

The uniqueness of the output corresponding to each input is described by the phrase “only one” in the definition. In addition, each element of the domain, A , is indeed an input and has an image in the range, B ; this a second, more technical property of functions. Relations that are not functions need not have this property. For example, in the order relation on the interval $[0, 1] \subset \mathbb{R}$, nothing is greater than 1, so 1 does not occur as the first coordinate of any pair (x, y) in the relation $\{(x, y) \in [0, 1] : x < y\}$. We can always arrange for the domain of a function to be restricted to elements that have outputs, but if we want to talk of an order relation *on* $[0, 1]$, we cannot leave 1 out!

Warning: as with many concepts, the terminology used about functions varies somewhat. Some authors use the word “range” to refer to what we have defined as the image, and use the word “codomain” to refer to what we have defined as the range. When reading any source, be sure to take note of how words are used in that context!

A misleading view of functions that is, unfortunately, rather commonplace in secondary schools is that they are “machines” that change inputs into outputs. Functions don’t change anything! The squaring function, for instance, associates the output 4 to the input 2, but it does not change 2 into 4. If you could find any means of changing 2 into 4, you would be doing better than the dreams of the alchemists who tried to

change lead into gold! Even the functions we speak of as “transformations,” such as linear transformations of the plane or of higher-dimensional vector spaces, do not change any of the positions or vectors, but only associate them to new positions or vectors. It is sometimes helpful to think of *representations* of the vectors as moving, but the vectors themselves do not change. Much is lost if one thinks of a function as somehow “producing” an output from an input, because it is the correspondence between inputs and outputs that characterizes a function. Seeing this correspondence (such as when one graphs the function) requires thinking of the input and output together, as an ordered pair.

6.8.2 Some special types of functions

Identity functions

Given a set A , the *identity function on A* is the function $i_A : A \rightarrow A$ defined by $i_A(a) = a$. (In set notation, $i_A = \{(x, y) \in A \times A : x = y\}$. Note that the identity function is the same as the equality relation on A , sometimes called the *diagonal* of $A \times A$.)

Characteristic functions

Let $B \subseteq A$. The *characteristic function of B* , denoted χ_B , is the function from A to $\{0, 1\}$ defined by

$$\chi_B(a) = \begin{cases} 1, & \text{if } a \in B \\ 0, & \text{if } a \notin B \end{cases}$$

The χ_B is sometimes called the *indicator function for B* ; it indicates whether or not an element of the domain is in B .

Binary operations

A *binary operation* on a set A is a function from $A \times A$ to A . This definition formalizes the notion introduced in the previous chapter.

Unitary operations

A *unitary operation* on a set A is a function from A to itself. For example, the function $f : \mathbb{R} \rightarrow \mathbb{R}$ defined by $f(x) = x^2$ is a unitary operation. (Recall that in set notation, $f = \{(x, y) \in \mathbb{R} \times \mathbb{R} : x^2 = y\}$.)

Taking additive and multiplicative inverses are important examples of unitary operations. Thus, the functions $f : \mathbb{R} \rightarrow \mathbb{R}$ defined by $f(x) = -x$ and $g : \mathbb{R} \setminus \{0\} \rightarrow \mathbb{R} \setminus \{0\}$ defined by $g(x) = \frac{1}{x}$ are unitary operations.

Observe that $f \circ f = i_{\mathbb{R}}$ and $g \circ g = i_{\mathbb{R} \setminus \{0\}}$. Any unitary operation whose composition with itself is the identity is called an *inversion*. Reflection in a line in the plane is a geometric example of an inversion, as is reflection through a point.

Metrics

A very important type of function that occurs in geometry is a *metric*, or “distance function.” A *metric* on a set A is a function $d : A \times A \rightarrow \mathbb{R}$ that satisfies the following properties:

- *Symmetry*: $\forall a, b \in A, d(a, b) = d(b, a)$. (Note that this is a different notion of symmetry from that of a relation *on* A being symmetric.)
- *Positive definiteness*: $\forall a, b \in A, d(a, b) \geq 0 \wedge (d(a, b) = 0 \Leftrightarrow a = b)$.
- *Triangle inequality*: $\forall a, b, c \in A, d(a, b) + d(b, c) \geq d(a, c)$.

6.8.3 Surjectivity, Injectivity, and Inverse Functions

Definition. A function is *surjective* if its image is equal to its range. That is, every element of the range is the image of some point in the domain. A technical formulation is more useful in proofs; $f : A \rightarrow B$ is *surjective* if, $\forall b \in B, \exists a \in A$ such that $f(a) = b$. (Note that the statement obtained by interchanging a and b in the quantifiers is part of the definition of a function: $\forall a \in A, \exists b \in B$ such that $f(a) = b$.) If $f : A \rightarrow B$ is surjective, we say it is a function from A *onto* B .

Notation. If f is a surjective function from A onto B , it is common to write $f : A \twoheadrightarrow B$.

Definition. A function is *injective* if each element of its range occurs as the image of at most one element of the domain. Again, a technical formulation is more useful in proofs; $f : A \rightarrow B$ is *injective* if, $\forall a, a' \in A, f(a) = f(a') \Rightarrow a = a'$. (Note that the converse implication is again part of the definition of a function.) If $f : A \rightarrow B$ is injective, we say it is a function from A *into* B .

Notation. If f is an injective function from A into B , it is common to write $f : A \hookrightarrow B$.

Definition. A function that is both injective and surjective is called *bijective*.

Definition. Let $g : A \rightarrow B$. A function f is a *left inverse (function)* for g if $f \circ g = i_A$.

Definition. Let $g : A \rightarrow B$. A function h is a *right inverse (function)* for g if $g \circ h = i_B$.

Theorem (Existence of Left or Right Inverse Functions). *A function has a left inverse if and only if it is injective. A function has a right inverse if and only if it is surjective. It follows that a function has both a left and a right inverse if and only if it is bijective.*

The proof of this theorem is contained in the exercises.

It is important to avoid confusion here. A right or left inverse *function* of a given function is not necessarily the same as the inverse *relation* of that function. If g is injective but not surjective, a left inverse f for g contains the inverse relation of g as a proper subset, because images under f must be chosen for those elements of the range of g that are not in the image of g . It does not matter how these images are chosen when defining f ; therefore, there is not a unique left inverse of such a function. If g is surjective but not injective, a right inverse h for g is a proper subset of the inverse relation of g , because an image under h for each element of the range of g must be chosen from among its pre-images. Again, this choice is not unique, so there is no unique right inverse of such a function. (To see all this more clearly, I suggest drawing diagrams. Working the exercises will also clarify these facts - as usual!) If g is bijective, then the inverse relation of g is a function and is both the unique left inverse and the unique right inverse of g . This function is called simply the *inverse (function)* of g and denoted by g^{-1} . (Please note again that this has nothing to do with reciprocals!) This result is stated below; its proof is left as an exercise.

Theorem (Existence of an Inverse Function). *Let $g : A \rightarrow B$ be a function. The inverse relation of g is a function if and only if g is bijective. In this case, the inverse relation of g is the unique left inverse function for g and also the unique right inverse function for g .*

Definition. Let $g : A \rightarrow B$ be bijective. The *inverse (function)* to g , denoted by g^{-1} , is the unique function such that $g^{-1} \circ g = i_A$ and $g \circ g^{-1} = i_B$.

Combining and summarizing the previous two theorems, we obtain:

Theorem. *A function g has both a left inverse f and a right inverse h if and only if g is bijective and $f = h = g^{-1}$.*

Definition. Given a function $f : A \rightarrow B$, for any $b \in B$ we define the *pre-image* of b , denoted by $f^{-1}(b)$, to be the set $\{a \in A : f(a) = b\}$. More generally, for $D \subseteq B$, the pre-image of D , denoted by $f^{-1}(D)$ is defined as the set $\{a \in A : f(a) \in D\}$. Note that the pre-image of a point or set may be empty.

When working with functions, it is important to be aware of context, as some symbols have different meanings in different contexts. For example, f^{-1} may be used to denote the inverse function, if it exists, or the pre-image of a point or set (even if the function f is not bijective and has no inverse function). If f is bijective, then f^{-1} generally means the inverse function; hence, if $f(a) = b$, $f^{-1}(b) = a$. In this case, the element a is called the pre-image of b , even though according to the above definition, the pre-image of b is the set $\{a\}$. Similarly, for any function $f : A \rightarrow B$ and any $b \in B$, $f^{-1}(b) = f^{-1}(\{b\})$. For obvious reasons, these multiple usages of words and notations do not generally cause confusion. That is why we use them! As is often the case, it would be unnecessarily cumbersome to introduce more complicated notation to avoid them. (Mathematicians sometimes refer to these multiple usages as “abuse of notation.” Better to abuse the notation than the reader!)

6.8.4 Exercises

1. Let $f : A \rightarrow B$ and $g : B \rightarrow C$ be functions. Prove that the composition $g \circ f$ is a function from A to C , and that $g \circ f(a) = g(f(a))$.

Remark. Some writers prefer to write functions on the right, instead of the left, so that compositions come out in a more natural order: $b = (a)f$, if $(a, b) \in f$. Then the composition above comes out as $f \circ g$, with $(a)f \circ g = ((a)f)g$. However nice, this practice has not generally caught on. Tradition dies hard!

2. Let $A = \{1, 2\}$ and $B = \{3, 4\}$. List all possible functions from A to B and sketch their graphs. Is it possible that a function from A to B is injective but not surjective? Is it possible that a function from A to B is surjective but not injective? Explain!
3. Let $A = \{1, 2\}$ and $B = \{3, 4\}$. Give the inverse function for each bijective function from A to B and sketch its graph. What is the geometric relationship between the graphs of a function and its inverse function?
4. Let $A = \{1, 2\}$ and $B = \{3, 4, 5\}$. List all possible functions from A to B and sketch their graphs. For each injective function, list all of its left inverses and sketch their graphs. Is it possible to give an example of a surjective function from A to B ? Explain.
5. Let $A = \{1, 2, 3\}$ and $B = \{3, 4\}$. List all possible functions from A to B and sketch their graphs. For each surjective function, list all of its right inverses. Is it possible to give an example of an injective function from A to B ? Explain.

6. Give an example of a set A and two functions, $f : A \rightarrow A$ and $g : A \rightarrow A$, such that $g \circ f \neq f \circ g$. (What does it mean for two functions to be equal?)
7. Let c be any real number. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by $f(x) = x + c$. Prove that f is bijective.
8. Let c be any real number other than zero. Let $g : \mathbb{R} \rightarrow \mathbb{R}$ be defined by $g(x) = cx$. Prove that g is bijective.
9. Let $h : \mathbb{R} \rightarrow \mathbb{R}$ be defined by $h(x) = x^2$. Prove that h is neither injective nor surjective.
10. Let $j : \mathbb{R} \rightarrow \mathbb{R}$ be defined by $j(x) = (x - 1)(x - 2)(x - 3)$. Sketch the graph of j (that is, the points $(x, y) \in \mathbb{R} \times \mathbb{R}$ such that $(x, y) \in j$, in other words, such that $y = j(x)$). Use the graph to present an informal argument that j is surjective but not injective.
11. With f, g, h , and j defined as above, give formulas for the following functions (you needn't simplify the complicated ones):
 - (a) $f \circ g$
 - (b) $g \circ f$
 - (c) $f \circ h$
 - (d) $h \circ f$
 - (e) $f \circ j$
 - (f) $j \circ f$
 - (g) $g \circ h$
 - (h) $h \circ g$
 - (i) $g \circ j$
 - (j) $j \circ g$
 - (k) $h \circ j$
 - (l) $j \circ h$
12. Prove: The composition of two injective functions is injective. The composition of two surjective functions is surjective. It follows that the composition of two bijective functions is bijective.

13. Let $f : A \rightarrow B$ and $g : B \rightarrow C$ be functions. Prove: The composition $g \circ f$ is injective only if f is injective. The composition $g \circ f$ is surjective only if g is surjective. It follows that the composition $g \circ f$ is bijective only if f is injective and g is surjective.
14. Show that the conditions given in the previous problem are not (nearly!) sufficient by giving an example of a pair of functions f and g such that f is injective and g is surjective, but $g \circ f$ is neither injective nor surjective.
15. Prove the theorem on the existence of left or right inverses.

Remark. The proof that a function has a right inverse if it is surjective requires that you make a choice from the pre-image of each element of the range, in order to define a right inverse. That you can make these choices all at once, even if the range is infinite, requires an additional set-theoretic assumption, called the Axiom of Choice. In essence, you are assuming the existence of a function, which you cannot define, that associates to each element of the range an element of its pre-image set. It is obviously necessary that each pre-image set be non-empty, but nothing in the nature of logic or in our previous axioms of set theory assures us that this condition is sufficient.

The Axiom of Choice is slightly controversial. Aside from invoking it in this exercise, we will not need or discuss the Axiom of Choice further in this book; however, you should develop an awareness of when you require such a *choice function*.

16. Prove the theorem on the existence of an inverse function.
17. Prove that for any relation R of A to B , $i_B \circ R = R = R \circ i_A$. In particular, for any function $f : A \rightarrow B$, $i_B \circ f = f = f \circ i_A$. Thus, the identity functions constitute a sort of identity for the composition operation.

6.8.5 The Induced Function on Power Sets

Any function $f : A \rightarrow B$ induces a function from the power set of A to the power set of B , which we generally denote by the same letter, since it is clear from context which function we mean. (Abuse of notation again!) This induced function is defined as follows: for each $C \subseteq A$, the image of C is defined as $f(C) = \{b \in B : (\exists c \in C)f(c) = b\}$. In other words, $f(C) = \{f(c) : c \in C\}$.

It is important to recognize the relationships between images and pre-images of sets, which the following exercises are intended to illuminate.

6.8.6 Exercises

In the following exercises, $f : A \rightarrow B$ is a function, $C \subseteq A$, $D \subseteq B$, $E \subseteq A$, and $F \subseteq B$.

1. Prove that $f(f^{-1}(D)) \subseteq D$. Give a counterexample to show that equality does not necessarily hold. Give a necessary and sufficient condition on f for equality to hold for all subsets $D \subseteq B$, and prove your condition is necessary and sufficient.
2. Prove that $f^{-1}(f(C)) \supseteq C$. Give a counterexample to show that equality does not necessarily hold. Give a necessary and sufficient condition on f for equality to hold for all subsets $C \subseteq A$, and prove your condition is necessary and sufficient.
3. Prove that $f(C \cup E) = f(C) \cup f(E)$.
4. Prove that $f(C \cap E) \subseteq f(C) \cap f(E)$. Give a counterexample that shows equality need not hold. Give a necessary and sufficient condition on f for equality to hold for all subsets C and E of A , and prove your condition is necessary and sufficient.
5. Prove that $f^{-1}(D \cup F) = f^{-1}(D) \cup f^{-1}(F)$.
6. Prove that $f^{-1}(D \cap F) = f^{-1}(D) \cap f^{-1}(F)$.

Remark. Observe that pre-images behave nicely under intersection as well as union, as shown in exercises 5 and 6, because each element of the domain has only one image. Images do not behave so nicely under intersection, because an element of the range may have multiple pre-images. (See exercise 4); However, this does not present a problem for unions, which behave nicely under images, too. (See exercise 3). Make a point of remembering these observations!

6.9 Properties of Relations: a Compendium

In working the examples, you may well have noticed relations that have special properties. Functions are an obvious example; order relations (which we will soon revisit in greater depth) are another. Remember that mathematics is about recognizing patterns, and recognizing patterns is about finding recurrent properties and naming them so we can talk about them, work with them, and deduce their consequences. Here we collect some distinctive general properties that a relation may - or may not - have.

- *Functionality:* A relation R of A to B is *functional* if, $\forall a \in A, \exists! b \in B$ such that $(a, b) \in R$.

As we have seen, a relation with this property is called a function. (Relations that are not functional still work very well in many situations; the mathematical usage of the term does not imply its colloquial meaning!)

The following properties of relations are generally applied to relations on a single set:

- *Comparability*: A relation R on A satisfies this property if any two distinct elements of A are comparable, that is, related in some way by R . Thus, R satisfies *comparability* if, $\forall a_1 \in A$ and $\forall a_2 \in A$, $a_1 \neq a_2 \Rightarrow (a_1, a_2) \in R \vee (a_2, a_1) \in R$.
- *Reflexivity*: A relation R on A is *reflexive* if, $\forall a \in A$, $(a, a) \in R$.
- *Anti-reflexivity*: A relation R on A is *anti-reflexive* if, $\forall a \in A$, $(a, a) \notin R$.
- *Symmetry*: A relation R on A is *symmetric* if, $\forall a_1 \in A$ and $\forall a_2 \in A$, $(a_1, a_2) \in R \Rightarrow (a_2, a_1) \in R$. (Note that since a_1 and a_2 are arbitrary, biconditionality is automatic.)
- *Anti-symmetry*: A relation R on A is *anti-symmetric* if, $\forall a_1 \in A$ and $\forall a_2 \in A$, $(a_1, a_2) \in R \Rightarrow (a_2, a_1) \notin R$. (Here biconditionality need *not* hold. If R satisfies comparability, the converse implication is of course true, by definition.)
- *Transitivity*: A relation R on A is *transitive* if, $\forall a_1 \in A$, $\forall a_2 \in A$, and $\forall a_3 \in A$, $[(a_1, a_2) \in R \wedge (a_2, a_3) \in R] \Rightarrow (a_1, a_3) \in R$.

As always, the warning applies that terminology varies somewhat among authors. Some authors use *irreflexive* instead of anti-reflexive, for example. Worse yet, some authors use *nonreflexive* to mean simply *not* reflexive, whereas others use it the way we have used anti-reflexive, which describes a much stronger condition. Pay careful attention to the definitions that pertain in any particular text!

6.9.1 Exercises

1. (Reformulation of transitivity.) Show that a relation R on A is transitive if and only if $R \circ R \subseteq R$.
2. Show that if $R \subseteq A \times B$ is any relation of A to B , then $R^{-1} \circ R$ is a symmetric relation on A and $R \circ R^{-1}$ is a symmetric relation on B .

6.10 Order Relations

Let R be a relation on a set A . As we know from studying the order of the real numbers, R is an *order relation* if it is anti-reflexive and transitive. As we have seen for the ordering of the real numbers and will prove generally in the exercises, anti-symmetry is an additional property of order relations; it is a consequence of transitivity together with anti-reflexivity. Usually the symbol $x < y$ or something similar such as $x \prec y$, if there a desire to emphatically distinguish R from the stand order of numbers, is used to denote the statement $(x, y) \in R$. (We may also use $x > y$ or $x \succ y$ if we want to think of the order in reverse. See Exercise ?? of Section 6.10.4.)

6.10.1 Weak Orders Versus Strong Orders

Some authors define an order relation slightly differently. They take order relations to be reflexive, and denote them using symbols such as \leq or \preceq . This necessitates modifying the anti-symmetry property, since for each element a of A , $(a, a) \in R$ obviously violates anti-symmetry. We require that this is the only situation in which anti-symmetry is violated: $\forall a_1 \in A$ and $\forall a_2 \in A$, $[(a_1, a_2) \in R \wedge (a_2, a_1) \in R] \Rightarrow a_1 = a_2$. (Equivalently, using the contrapositive: $a_1 \neq a_2 \Rightarrow [(a_1, a_2) \in R \Rightarrow (a_2, a_1) \notin R]$). Let us call this property *weak anti-symmetry*. An order relation in either conception must be transitive.

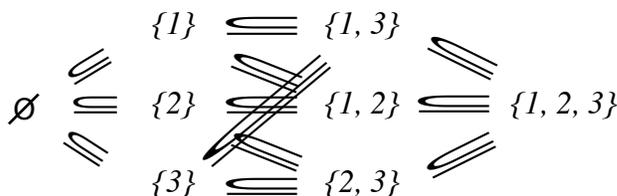
To avoid ambiguity, an order relation as we originally defined it (anti-reflexive, transitive, and consequently anti-symmetric) is called a *strong* order, and an order relation defined in the second manner (reflexive, weakly anti-symmetric, transitive) is called a *weak* order. If we use the word “order” unmodified, we mean a strong one. (Some authors use the phrase “strict order” for a strong order and, as noted above, take the unmodified word “order” to mean a weak one. The meaning is usually clear from notation and context in any case.)

This is all a matter of convention, not substance. If $<$ is a strong order on a set A , then defining $a_1 \leq a_2$ by $a_1 < a_2 \vee a_1 = a_2$ gives a weak order. Conversely, if \leq is a weak order on A , then defining $a_1 < a_2$ by $a_1 \leq a_2 \wedge a_1 \neq a_2$ gives a strong order. The proof of these fairly obvious facts is left to the exercises. Depending on the situation, it may be more natural to define a weak order first and the corresponding strong order in terms of it, or it may be more natural to define a strong order first and the corresponding weak order in terms of it.

6.10.2 Partial Orders Versus Total Orders

An order relation is *total* if it satisfies comparability. Otherwise, it is called *partial*. Visually, one can think of a total order on A as arranging the elements of A in a line, whereas a partial order arranges them in a sort of directed web.

Example. Let $A = \{1, 2, 3\}$. The partial ordering of $\mathcal{P}(A)$ by inclusion is illustrated in the following figure:



6.10.3 The Dictionary Order on a Cartesian Product

Let $<_A$ be a (strong) order on a set A , and let $<_B$ be a strong order on a set B . Then we define a (strong) order on $A \times B$ by:

$$(a, b) < (a', b') \Leftrightarrow [a <_A a'] \vee [a = a' \wedge b <_B b'].$$

The order defined in this manner is called the *dictionary order* on $A \times B$ (corresponding to the given orders $<_A$ and $<_B$), for reasons that should be obvious.

Example. Consider the standard order on \mathbb{R} . In the corresponding dictionary order on $\mathbb{R} \times \mathbb{R}$, $(0, 2) < (1, 0) < (1, 1)$. Note: do not confuse the *interval* (x, y) , defined as $\{z \in \mathbb{R} : x < z < y\}$, with the *ordered pair* (x, y) . Although the same notation is used, the meaning should be clear from context.

6.10.4 Exercises

1. Prove that if a relation on a set A is anti-reflexive and transitive, then it is anti-symmetric. (Hint: The proof that worked in for the standard order on the real numbers works in general. See if you can remember it without looking!)

2. Let \leq be a weak order relation on a set A . (That is, assume \leq is reflexive, weakly anti-symmetric, and transitive.) Define $a_1 < a_2$ to mean that $a_1 \leq a_2$ and $a_1 \neq a_2$. Prove that $<$ is anti-reflexive and transitive, hence a strong order relation. (Anti-reflexivity is immediate from the definition; transitivity takes a bit more care.)
3. Let $<$ be a strong order relation on a set A . (That is, assume $<$ is anti-reflexive and transitive, hence also anti-symmetric.) Define $a_1 \leq a_2$ to mean that $a_1 < a_2$ or $a_1 = a_2$. Prove that \leq is reflexive, weakly anti-symmetric, and transitive, hence a weak order relation. (Reflexivity is immediate from the definition.)
4. Demonstrate a relation on the two-element set $\{a, b\}$ ($a \neq b$) that is reflexive and transitive but not weakly anti-symmetric. This shows that the condition of weak anti-symmetry must be assumed for a weak order relation. (Hint: This is easy; you don't have a whole lot of choice, here!)
5. If $<$ is a total order on a set A , it follows immediately from comparability that, for any elements $a_1, a_2 \in A$, one of the following holds:

- $a_1 = a_2$,
- $a_1 < a_2$, or
- $a_2 < a_1$.

(Make sure you understand why!) Complete the proof of trichotomy in this general context by proving that *only one* of these properties holds. (Hint: The proof that worked for the standard order on the real numbers works in general. See if you can remember it without looking!)

6. Let $<$ be an order on a set A . Define $a_1 > a_2$ in the usual manner to mean $a_2 < a_1$. (Thus the relation $>$ is just the inverse of the relation $<$.) Show that $>$ is also an order relation on A . (That is, show that $>$ is anti-reflexive and transitive.)
7. Verify that the dictionary order is indeed anti-reflexive and transitive.
8. Define a relation on $\mathbb{R} \times \mathbb{R}$ as follows: $(x_1, y_1) \prec (x_2, y_2) \Leftrightarrow y_1^2 - x_1 < y_2^2 - x_2 \vee (y_1^2 - x_1 = y_2^2 - x_2 \wedge y_1 < y_2)$. Show that \prec is an order relation.
9. Draw a diagram depicting the partial order by inclusion on $\mathcal{P}(\{1, 2, 3, 4\})$.

6.10.5 Maxima and Minima; Bounds; Suprema and Infima

Given an order relation on a set, we often want to compare some element to all of the elements in a specified subset. A given subset may or may not contain an element that is greater than all the others. If it does, this element is unique. The reasoning is so simple that the result takes longer to state than to prove!

Proposition. *Let $<$ be a (strong) order on a set A (with corresponding weak order \leq). Let $B \subseteq A$. Suppose $b_0 \in B$ has the property that $(\forall b \in B)b \leq b_0$. Suppose $b'_0 \in B$ also has this property. Then $b_0 = b'_0$.*

Proof. By assumption, $b_0 \leq b'_0$ and $b'_0 \leq b_0$. Thus $b_0 = b'_0$ (by anti-symmetry). \square

We are therefore justified in using the word *the* in the following definition:

Definition. Let $<$ be a (strong) order on a set A (with corresponding weak order \leq). Let $B \subseteq A$. The element $b_0 \in B$ is the *greatest element* of B , also called the *maximum* of B , if $(\forall b \in B)b \leq b_0$.

Similarly (and as we have already seen in the section on well-ordering):

Definition. Let $<$ be a (strong) order on a set A (with corresponding weak order \leq). Let $B \subseteq A$. The element $b_0 \in B$ is the *smallest element* of B , also called the *least element* or *minimum* of B , if $(\forall b \in B)b_0 \leq b$.

Example. Consider the standard order relation on the real numbers. The interval $(0, 1)$ has no greatest or least element. (Recall that the *interval* $(0, 1)$ is the set $\{x \in \mathbb{R} : 0 < x < 1\}$.)

Example. Again considering the standard order on \mathbb{R} , 1 is the greatest element of the interval $(0, 1]$. (Recall that $(0, 1] = \{x \in \mathbb{R} : 0 < x \leq 1\}$.) The interval $(0, 1]$ has no least element.

If an element is at least as large as all of the elements in a given subset, it is called an *upper bound* for that subset; an upper bound need not be an element of the given subset. Similarly, an element that is at least as small as every element in a given subset is called a *lower bound* for that subset. Formally:

Definition. Let $<$ be a (strong) order on a set A , (with corresponding weak order denoted by \leq). Let $B \subseteq A$. An element $a \in A$ is an *upper bound* for B if, $\forall b \in B, b \leq a$.

Definition. Let $<$ be a (strong) order on a set A , (with corresponding weak order denoted by \leq). Let $B \subseteq A$. An element $a \in A$ is an *lower bound* for B if, $\forall b \in B, a \leq b$.

Upper and lower bounds are generally not unique; anything larger than an upper bound is still an upper bound, and anything smaller than a lower bound is still a lower bound. Note that, for the empty set, any element is both an upper bound and a lower bound, since the defining conditions are vacuously true. Also note that an upper or lower bound for a given set need not exist.

A set that has an upper bound is said to be *bounded above*; a set that has a lower bound is said to be *bounded below*. A set that is bounded both above and below is said to be *bounded*.

Example. Consider the standard order relation on the real numbers. The number 1 is an upper bound for the interval $(0, 1)$. So is the number 2, or any number larger than 1. Any number less than or equal to 0 is a lower bound for $(0, 1)$.

Example. Again considering the standard order on \mathbb{R} , 1 is an upper bound for the interval $(0, 1]$, as is any number greater than 1. Any number less than or equal to 0 is a lower bound for $(0, 1]$.

Example. Again considering the standard order on \mathbb{R} , 0 is a lower bound for the interval $(0, \infty) = \{x \in \mathbb{R} : 0 < x\}$, as is any number less than 0. The interval $(0, \infty)$ has no upper bound.

Example. In the dictionary order on $\mathbb{R} \times \mathbb{R}$, $(1, 0)$ is an upper bound for the set $\{0\} \times \mathbb{R} = \{(x, y) \in \mathbb{R} \times \mathbb{R} : x = 0\}$. So are $(\frac{1}{2}, 0)$, $(\frac{1}{4}, 0)$, $(\frac{1}{8}, 0)$, $(\frac{1}{8}, 1)$, and any ordered pair (x, y) such that $x > 0$. Similarly, any ordered pair (x, y) such that $x < 0$ is a lower bound for $\{0\} \times \mathbb{R}$.

Example. In the dictionary order on $\mathbb{R} \times \mathbb{R}$, the set $\mathbb{R} \times \{0\} = \{(x, y) \in \mathbb{R} \times \mathbb{R} : y = 0\}$ has neither an upper nor a lower bound.

If a set is bounded above, the set of upper bounds may or may not have a least element. If it does, that least element is unique, as we proved above, so we are justified in using the word *the* in the following definition:

Definition. Let $<$ be an order on a set A , and let $B \subseteq A$. The element a is the *least upper bound*, also called the *supremum*, for B if:

- a is an upper bound for B , and
- if c is any upper bound for B , the $a \leq c$.

In other words, a is the least upper bound (supremum) for B if it is the least element of the set of upper bounds for B .

Similarly:

Definition. Let $<$ be an order on a set A , and let $B \subseteq A$. An element a is the *greatest lower bound*, also called the *infimum*, for B if:

- a is a lower bound for B , and
- if c is any lower bound for B , the $a \geq c$.

In other words, a is the greatest lower bound (infimum) for B if it is the greatest element of the set of lower bounds for B .

Again, note that a given set need have a supremum or infimum, even if it is bounded above or below. Also note that the supremum or infimum, if it exists, need not be an element of the set.

Notation. The supremum for a set $B \subseteq A$ is denoted by $\sup_A B$. The infimum for a set B is denoted by $\inf_A B$. If the ordered set A is clear from context, it is generally omitted. If there is more than one ordering of A being considered, the order to which the supremum or infimum refers may also be included in the subscript, as in $\inf_{A, <} B$.

“Supremum” and “infimum” are Latin words; therefore, the plural of “supremum” is “suprema,” and the plural of “infimum” is “infima.”

Examples. $\sup_{\mathbb{R}}(0, 1) = 1$; $\sup_{\mathbb{R}}(0, 1] = 1$; $\inf_{\mathbb{R}}(0, 1) = 0$; $\inf_{\mathbb{R}}(0, 1] = 0$; $\inf_{\mathbb{R}}(0, \infty) = 0$. $(0, \infty)$ has no supremum, because it is not even bounded above. Similarly, $\mathbb{R} \times \{0\}$ has neither a supremum nor an infimum in the dictionary order on $\mathbb{R} \times \mathbb{R}$. Even though $\{0\} \times \mathbb{R}$ is bounded (above and below) in the dictionary order on $\mathbb{R} \times \mathbb{R}$, it has neither a supremum nor an infimum. There is no smallest element (x, y) with $x > 0$; for any $x > 0$, there is always a number x' such that $x > x' > 0$. Similarly, there is no largest element (x, y) with $x < 0$.

Finally, an element of a set is called *maximal* if no element of the given set is greater. For a total order, the concepts of maximal and maximum are the same, but for a partial order they are not. A maximal element need not be greater than every other element of the given set; it may simply be incomparable to some of them. There may be more than one maximal element. However, if a partially ordered set has a maximum, the maximum is the unique maximal element. Similarly, an element of a set is called *minimal* if no element of the given set is smaller.

Example. Consider the partial order on $\mathcal{P}(\{1, 2, 3\})$ given by inclusion. The set $\{1, 2, 3\}$ is both maximal and maximum. The empty set is both minimal and minimum.

Example. If we consider the family $\mathcal{P}(\{1, 2, 3\}) \setminus \{\{1, 2, 3\}\}$ of *proper* subsets of $\{1, 2, 3\}$, then each of the sets $\{1, 2\}$, $\{1, 3\}$, and $\{2, 3\}$ is maximal. There is no maximum.

6.10.6 Exercises

1. Verify that the least element of a set is unique.
2. Prove that if a set has a maximum, the maximum is also the supremum for that set.
3. Prove that if a set has a minimum, the minimum is also the infimum for that set.
4. Prove that a set has a maximum if and only if an element of the set is an upper bound for the set.
5. Prove that a set has a minimum if and only if an element of the set is a lower bound for the set.
6. Prove that if a set has a maximum, the maximum is the unique maximal element of that set.
7. Prove that if a set has a minimum, the minimum is the unique minimal element of that set.

In the following exercises, \mathbb{R} is assumed to have the standard order, and $\mathbb{R} \times \mathbb{R}$ is assumed to have the corresponding dictionary order.

8. For each of the following sets, state if it has a maximum and, if so, state what the maximum is.
 - (a) The interval $(-\infty, 2)$.
 - (b) The interval $(-\infty, 2]$.
 - (c) $\{4, 5, 6, 7\}$.
 - (d) $\{\frac{1}{2}, \frac{2}{3}, \frac{3}{4}, \dots\}$.
 - (e) $\{\frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \dots\}$.
 - (f) \mathbb{N} .
 - (g) $\mathbb{N} \times \{2\}$.
 - (h) $\{2\} \times \mathbb{N}$.
 - (i) $(0, 1) \times \mathbb{N}$.
 - (j) $[0, 1] \times \mathbb{N}$.

9. For each of the following sets, state if it is bounded above. If it is, state if it has a least upper bound, and if so, state what the least upper bound is.
- (a) The interval $(-\infty, 2)$.
 - (b) The interval $(-\infty, 2]$.
 - (c) $\{4, 5, 6, 7\}$.
 - (d) $\{\frac{1}{2}, \frac{2}{3}, \frac{3}{4}, \dots\}$.
 - (e) $\{\frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \dots\}$.
 - (f) \mathbb{N} .
 - (g) $\mathbb{N} \times \{2\}$.
 - (h) $\{2\} \times \mathbb{N}$.
 - (i) $(0, 1) \times \mathbb{N}$.
 - (j) $[0, 1] \times \mathbb{N}$.
10. For each of the following sets, state if it has a minimum and, if so, state what the minimum is.
- (a) The interval $(-\infty, 2)$.
 - (b) The interval $(-\infty, 2]$.
 - (c) $\{4, 5, 6, 7\}$.
 - (d) $\{\frac{1}{2}, \frac{2}{3}, \frac{3}{4}, \dots\}$.
 - (e) $\{\frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \dots\}$.
 - (f) \mathbb{N} .
 - (g) $\mathbb{N} \times \{2\}$.
 - (h) $\{2\} \times \mathbb{N}$.
 - (i) $(0, 1) \times \mathbb{N}$.
 - (j) $[0, 1] \times \mathbb{N}$.
11. For each of the following sets, state if it is bounded below. If it is, state if it has a greatest lower bound, and if so, state what the greatest lower bound is.
- (a) The interval $(-\infty, 2)$.

- (b) The interval $(-\infty, 2]$.
- (c) $\{4, 5, 6, 7\}$.
- (d) $\{\frac{1}{2}, \frac{2}{3}, \frac{3}{4}, \dots\}$.
- (e) $\{\frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \dots\}$.
- (f) \mathbb{N} .
- (g) $\mathbb{N} \times \{2\}$.
- (h) $\{2\} \times \mathbb{N}$.
- (i) $(0, 1) \times \mathbb{N}$.
- (j) $[0, 1] \times \mathbb{N}$.

For the following exercises, consider $\mathbb{R} \times \mathbb{R}$ with the order given in Exercise 8 of Section 6.10.4.

- 12. For each of the following sets, state if it has a maximum and, if so, state what the maximum is.
 - (a) $\mathbb{R} \times \{2\}$.
 - (b) $\{2\} \times \mathbb{R}$.
- 13. For each of the following sets, state if it is bounded above. If it is, state if it has a least upper bound, and if so, state what the least upper bound is.
 - (a) $\mathbb{R} \times \{2\}$.
 - (b) $\{2\} \times \mathbb{R}$.
- 14. For each of the following sets, state if it has a minimum and, if so, state what the minimum is.
 - (a) $\mathbb{R} \times \{2\}$.
 - (b) $\{2\} \times \mathbb{R}$.
- 15. For each of the following sets, state if it is bounded below. If it is, state if it has a greatest lower bound, and if so, state what the greatest lower bound is.
 - (a) $\mathbb{R} \times \{2\}$.

(b) $\{2\} \times \mathbb{R}$.

16. Let $A = \{1, 2, 3, 4\}$. What are the maximal and minimal elements of $\mathcal{P}(A)$, ordered by inclusion?
17. Let $A = \{1, 2, 3, 4\}$, and let $B = \mathcal{P}(A) \setminus \{A, \emptyset\}$, ordered by inclusion. What are the maximal and minimal elements of B ?

6.11 Equivalence Relations

As the name suggests, equivalence relations are used to group elements that have some common quality, ignoring the differences among them in other respects. As an example, suppose we have a metric that measures distances in the plane. Then we can define two segments to be *congruent* if they have the same length. (That is, if the distances between their respective endpoints are the same.) In doing so, we ignore the various positions of the segments and focus only on length as the quality of interest. Congruence is an equivalence relation.

In general, an *equivalence relation* is defined to be any relation that is reflexive, symmetric, and transitive. We can readily check these properties for congruence of segments in the plane: a segment obviously has the same length as itself; if segment AB has the same length as CD , then obviously CD has the same length as AB ; if AB has the same length as CD , and CD has the same length as EF , then clearly AB has the same length as EF . “Sameness” is, by its very nature, reflexive, symmetric, and transitive! Conversely, any relation with these three properties, even if it is defined in a manner that does not make them quite so obvious, defines a way in which related objects are the same. If R is an equivalence relation on a set A , we generally use notation such as $a_1 \sim a_2$, $a_1 \simeq a_2$, $a_1 \approx a_2$, $a_1 \cong a_2$, or $a_1 \equiv a_2$ to denote that $(a_1, a_2) \in R$.

Congruence is a special case of a very general method of defining an equivalence relation. Let $f : A \rightarrow B$ be any function. Define $a_1 \simeq_f a_2$ to mean $f(a_1) = f(a_2)$. Then \simeq_f is an equivalence relation. (For the congruence relation on segments, segments correspond to pairs of points. If d denotes the metric, define $f(AB) = d(A, B)$. Then \simeq_f is the congruence relation on segments.)

6.11.1 Equivalence Classes and Partitions

Let \simeq be an equivalence relation on a set A . Let a be an element of A . We define the *equivalence class* of a , denoted by $[a]$, to be the set of all elements of A that are equivalent to a : $[a] = \{a' \in A : a' \simeq a\}$. (If there might be ambiguity about which particular

equivalence relation on A is being considered, it may be specified as a subscript, as in $[a]_{\simeq}$. Understanding and describing the equivalence classes provides a good way to “picture” an equivalence relation.

Example. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by $f(x) = |x|$. Then $(\forall x \in \mathbb{R})[x]_{\simeq_f} = \{\pm x\}$. Thus, the equivalence class of 0 contains a single element (0 itself), and every other equivalence class contains two elements. Note that the equivalence classes of a number and its additive inverse are the same: $[x] = [-x]$. If we graph \simeq_f in the Cartesian plane in the usual manner, we can see that each vertical line intersects the graph at two points, except for the y -axis.

It should be intuitively clear that the equivalence classes divide A into disjoint subsets. Each subset is a “largest possible” set of elements that are equivalent, in the sense that nothing else could be added, because nothing else is equivalent to those elements already in the equivalence class. We now make this notion precise and prove it.

Definition. Let A be a set. A *partition* of A is a family $\mathcal{B} \subset \mathcal{P}(A)$ of non-empty subsets of A such that:

- Every element of A is in some subset in the family: $(\forall a \in A)(\exists B \in \mathcal{B})a \in B$. Equivalently (by definition), $\bigcup_{B \in \mathcal{B}} B = A$.
- The distinct subsets of the family are pairwise disjoint: $(\forall B_1, B_2 \in \mathcal{B})(B_1 \neq B_2 \Rightarrow B_1 \cap B_2 = \emptyset)$.

In other words, together the criteria for a partition specify that every element of A is in exactly one subset of the partition, as the term suggests. (The first criterion specifies that every element is in at least one, and the second specifies that every element is in at most one.)

Theorem. Let \simeq be an equivalence relation on a set A . Let $\mathcal{B} = \{[a] : a \in A\}$. The family \mathcal{B} is a partition of A .

Proof.

Claim. The sets in \mathcal{B} are non-empty, and every element of A is in one of them.

Let $a \in A$. By reflexivity, $a \in [a]$. (That is, every element is in its own equivalence class. Since each set in \mathcal{B} is, by definition, an equivalence class of some element of A , it is not empty.)

Claim. The sets in \mathcal{B} are pairwise disjoint: $(\forall B_1, B_2 \in \mathcal{B})(B_1 \neq B_2 \Rightarrow B_1 \cap B_2 = \emptyset)$.

Equivalently, $B_1 \cap B_2 \neq \emptyset \Rightarrow B_1 = B_2$. Suppose $[a_1] \cap [a_2] \neq \emptyset$. Then there is some element $a' \in [a_1] \cap [a_2]$; by definition of $[a_1]$, $[a_2]$, and intersection, $a' \simeq a_1$ and $a' \simeq a_2$. By transitivity and symmetry, we have that $a \in [a_1] \Leftrightarrow a \simeq a_1 \simeq a' \simeq a_2 \Leftrightarrow a \in [a_2]$. So $[a_1] = [a_2]$.

□

6.11.2 Exercises

1. Let $f : A \rightarrow B$. Verify that \simeq_f is an equivalence relation on A .
2. Consider $f : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ given by $f(x, y) = x + y$. What do the equivalence classes of \simeq_f look like?
3. Consider $g : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ given by $g(x, y) = x^2 + y^2$. What do the equivalence classes of \simeq_g look like?