

[tentative] Introduction to Abstract Mathematics  
through Inquiry  
M325K

Brian Katz  
Michael Starbird

January 28, 2009



# Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Developing Mathematical Ideas . . . . .	5
<b>2</b>	<b>Graph Theory</b>	<b>7</b>
2.1	The Königsberg Bridge Problem . . . . .	7
2.2	Connections . . . . .	8
2.3	Taking a Walk . . . . .	16
2.4	Trees . . . . .	23
2.5	Planarity . . . . .	24
2.6	Euler Characteristic . . . . .	26
2.7	Colorability . . . . .	29
2.8	Regular Planar Graphs . . . . .	34
2.9	Morphisms . . . . .	35
<b>3</b>	<b>Group Theory</b>	<b>37</b>
3.1	Examples Lead to Concepts . . . . .	37
3.2	Symmetry Groups of Regular Polygons . . . . .	45
3.3	Subgroups, Generators, and Cyclic Groups . . . . .	46
3.4	Products of Groups . . . . .	51
3.5	Symmetric Groups . . . . .	52
3.6	Maps between Groups . . . . .	54
3.7	Sizes of Subgroups and Orders of Elements . . . . .	60
3.8	Normal Subgroups . . . . .	61
3.9	Quotient Groups . . . . .	62
3.10	More Examples . . . . .	64
3.11	Groups in Action . . . . .	65
3.12	The Man Behind the Curtain . . . . .	68

<b>4</b>	<b>Calculus</b>	<b>71</b>
4.1	Know your Whereabouts . . . . .	71
4.2	Convergence . . . . .	72
4.3	Finding Limits . . . . .	81
4.4	Continuity . . . . .	91
4.5	Speeding and Zeno's Paradox <sup>TM</sup> . . . . .	98
4.6	Derivatives . . . . .	102
4.7	Speedometer Movie and Position . . . . .	107
4.8	Fundamental Theorem of Calculus . . . . .	108
<b>5</b>	<b>Topology - Math from Math</b>	<b>111</b>
5.1	Closeness . . . . .	112
5.2	Definition of a Topology . . . . .	113
5.3	Closed Sets . . . . .	118
5.4	Subspaces . . . . .	122
5.5	Bases . . . . .	122
5.6	Product Topologies . . . . .	124
5.7	Maps between Topological Spaces - Continuity . . . . .	125
5.8	Reasons behind the Madness . . . . .	127
5.9	Connected Spaces . . . . .	128
5.10	Metric Topologies and Continuity . . . . .	129
5.11	Abstraction is Useful . . . . .	131
<b>A</b>	<b>Appendix: Sets and Functions</b>	<b>133</b>
A.1	Sets . . . . .	133
A.2	Functions . . . . .	136
A.3	Special Functions . . . . .	138
A.4	Binary Operations . . . . .	140
A.5	Cardinality . . . . .	141
A.6	Notation . . . . .	142

# Chapter 1

## Introduction

### 1.1 Developing Mathematical Ideas

All mathematical ideas originate from human experience. We took our first shaky steps toward abstract mathematics when as toddlers we learned to count. Three cars, three bananas, and three dogs are physical realities that we can see and touch, but ‘three’ is not a concrete thing. The counting numbers are associated with collections of actual physical objects, but the counting numbers themselves give us our first abstract mathematical structure.

We soon learn to add numbers, multiply them, factor them, compare them, and otherwise discover and explore patterns, operations, and relationships among numbers. Numbers and their rich properties illustrate a strategy of creating and exploring concepts by starting with real world experiences and isolating features that then become mathematical ideas.

When we focus on the idea of measuring quantity in the world, we naturally develop mathematical concepts of number. When we focus on our visual or tactile impressions of the world, we develop geometrical ideas that range from Euclidean geometry to topology. When we isolate ideas of connections, we develop ideas of graph theory. When we analyze patterns and transformations, we find structures that lead to group theory. When we focus on change and motion, we are led to ideas of calculus.

Once a mathematical concept has begun its life as an abstraction of reality, then it takes on a reality of its own. We find variations and abstractions of ideas. For example, abstract extensions of the counting numbers include negative numbers, real numbers, and complex numbers. And the relationships and ways of combining counting numbers are extended, varied,

or abstracted to accommodate these new classes of numbers. Similarly, every mathematical concept propagates an ever-growing family of extensions, variations, and abstractions.

This book strives to demonstrate some of the basic strategies through which mathematical structures and concepts are created and developed. We treat graph theory, group theory, calculus, and topology in turn, showing how ideas are developed in each of these mathematical areas, but also demonstrating the commonalities in how abstract mathematics is discovered and explored.

One of the most basic features of mathematics is that human beings create it or discover it. Exploring mathematical ideas is an active process. You will not understand mathematical thought unless you personally participate in mathematical investigations. So this book actually is an invitation to you to think through the development of various mathematical concepts with the aid of our guidance. The fun of mathematics is to do it yourself. We have tried to design the experience to maximize the satisfaction you will feel in making mathematical ideas your own.

This book fundamentally consists of a series of exercises and theorem statements designed to introduce the “reader” to mathematical thought. We put “reader” in quotes because reading couldn’t be farther from your role. The most important part of the text is the part that isn’t there—the part you provide. The text primarily presents you with a series of challenges. In meeting those challenges by answering the questions, playing with examples, and proving the theorems on your own, you will develop intuition about particular mathematical concepts. You will also develop skills in how to investigate mathematical ideas and how to prove theorems on your own. This book strives to help you see the wonder of mathematical exploration. We hope you enjoy the journey.

## Chapter 2

# Graph Theory

### 2.1 The Königsberg Bridge Problem

Turn back the clock to the early 1700's and imagine yourself in the city of Königsberg, East Prussia. Königsberg was nestled on an island and on the surrounding banks at the confluence of two rivers. Seven bridges spanned the rivers as pictured below.



One day, Königsberg resident Friedrich ran into his friend Otto at the local Sternbuck's coffee shop. Otto bet Friedrich a Venti Raspberry Mocha Cappuccino that Friedrich could not leave the café, walk over all seven bridges without crossing over the same bridge twice (without swimming or flying), and return to the café. Friedrich set out, but never returned.

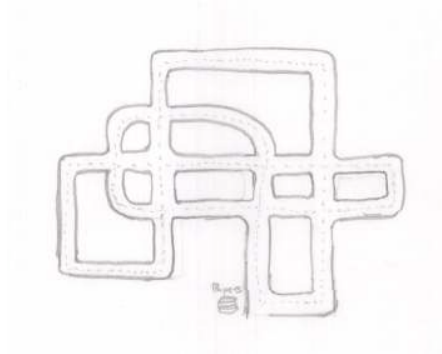
The problem of whether it is possible to walk over all seven bridges without crossing over the same bridge twice became known as the *Königsberg Bridge Problem*. As far as we know, Friedrich is still traipsing around the bridges of Königsberg, but a mathematician named Leonhard Euler did solve the Königsberg Bridge Problem in 1736, and his solution led to the modern area of mathematics known as *graph theory*.

## 2.2 Connections

One of the richest sources for developing mathematical ideas is to start with one or more specific problems and pare them down to their essentials. As we isolate the essential issues in specific problems, we create techniques and concepts that often have much wider applicability.

Sometimes it's quite hard to isolate the essential information from a single problem. If we consider several problems that “feel” similar, often the feeling of similarity guides us to the essential ingredients. It's a little like how, when playing games like *Catch Phrase* or *Taboo*, you choose several other words that have the secret word as a common thread. This process is important in creating the subject of *graph theory*. So let's begin by considering several additional questions that feel similar to the Königsberg Bridge Problem.

*The Paperperson's Puzzle.* Flipper, the paperperson, has a paper route in a residential area. Each morning at 5:00 a.m. a pile of papers is delivered to a corner in her neighborhood pictured below.



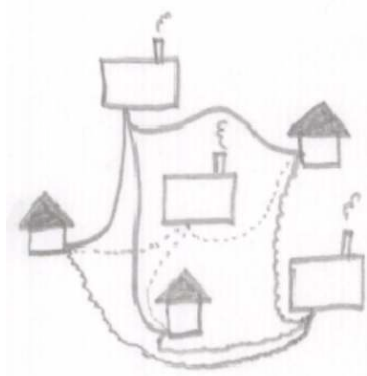
She puts all the papers in the basket of her bike and rides around the neighborhood flipping the papers in the general direction of the subscribers' houses. She rides down the middle of the streets and throws papers on both sides as she goes. When she finishes her route, she returns the leftover papers to the same location that she started. The question is whether she can complete her route without having to ride over the same street more than once.

The Königsberg Bridge Problem and the Paperperson's Puzzle have the similarity of taking a journey and returning to the starting point. However, some additional questions have similarities even though they do not involve



motion.

*The Gas-Water-Electricity Dilemma.* Three new houses have just been built in Houseville, and they all need natural gas, water, and electricity lines, each of which is supplied by a different company as pictured below.



Can each utility company lay a utility line to each house without having any of the utility lines cross?

*The Five Station Quandry.* Casey Jones wanted to build an elaborate model train set. He set up five stations and wanted to run tracks that connected each station to every other station. Could he build his layout with no crossing tracks or bridges?

Before you try to solve these problems, stop for a moment. What features of these problems are similar to one another? *Do not go on* until you think of at least one similarity among the problems.

Story problems are the bane of existence to non-mathematically oriented people, but mathematicians know exactly how to begin, namely, *abstraction*; that is, to isolate the salient information and to ignore the irrelevant information. The abstract concepts and techniques we create will not only help us solve these problems but will also be applicable to any other problem whose abstract essence is the same.

Here we will discuss the strategy of abstraction in the context of the Königsberg Bridge Problem, but please take analogous steps for the other five problems as well.

In the Königsberg Bridge Problem, what is important about the picture of the city? Does it matter how big the island is? Does it matter how long any of the bridges are? Does it matter that there are two bridges between

the northwest sector of town and the island? Ask yourself, “Which features of the problem set-up are relevant, and which features are not?” Asking yourself these questions is a big step towards mathematical maturity, and helping you to adopt the habit of asking yourself effective questions is one of the major goals of this book.

The features that seem to matter for the Königsberg Bridge Problem are the different locations (three land masses and the island) and the different bridges that cross between pairs of those locations. So one way to abstract the essence of this situation is to draw a dot for each location and a little line segment or edge for each bridge that connects a pair of dots (land masses). Since the problem does not ask about distance, the abstraction need not attempt to reflect any of the distances involved. Similarly, the physical locations of the land masses does not affect the problem, so the dots do not need to be positioned in any way that reflects the original city layout. The essential ingredients are locations and connections.

*Exercise 2.1.* Draw an abstracted picture that corresponds to the Königsberg Bridge Problem. Your picture should consist of dots and lines. Explain in your own words why this is a good representation.

*Exercise 2.2.* Draw similar abstracted pictures for each of the other challenges described above, the Paperperson’s Puzzle, the Gas-Water-Electricity Dilemma, and the Five Station Quandry. Your pictures should consist of dots and lines. Explain in your own words why these are good representations. In each case, what do your dots represent? What do your lines represent?

In some cases above, we must face the issue that we may not be able to draw all the lines connecting the dots without having the lines cross on the page. You will need to devise a strategy for indicating when an intersection of lines in your representation really shouldn’t be there. There are several standard ways to “fix” the drawings to remove the ambiguity.

*Exercise 2.3.* Attempt to solve the Five Station Quandry. If you cannot, draw an arrangement of the stations with railroad crossings using dots and lines where you create a notation for the unintentional crossings.

Perhaps you can think of alternative ways to abstract the essence of the situation that does not use a picture at all and hence avoids this issue altogether.

*Exercise 2.4.* Describe a new system for representing these situations that does not involve dots and lines but still contains the same information about the connections. Represent the data of one of the challenges above in your new system.

All of these problems resulted in abstractions that have similar characteristics. The visual representations all had dots and lines where each line connected two dots. Your non-visual representations probably used letters to represent locations or houses and utilities or people, while connections between pairs were indicated somehow, perhaps by writing down the pairs of letters that had a connection. Both the visual representation and the written one contained the same information and that information has two basic ingredients—some things and some connections between pairs of things. Once we have isolated these ingredients, we are ready to take an important step in the development of our concept, and that is to *make some definitions*.

Notice that we didn't start with the definitions. This process is typical of mathematical invention: we explore one or more situations that contain some intuitive or vague ideas in common and then we pin down those ideas by making a formal definition. Definitions are a mathematician's life's blood because they allow us to be completely clear about what is important and what is not important in a statement.

In all our examples above, including the Königsberg bridges, the train tracks, and the utilities and houses, we isolated the important features as things and connections. So we are ready to make a definition that captures situations of that type. The word we use to capture this abstract situation is a *graph*. Finally, here is our definition.

*Definition.* Let  $V$  be a finite, nonempty set, and let  $E$  be a set consisting of pairs of elements of the form  $\{v, v'\}$ , where  $v$  and  $v'$  are in  $V$ . Then the pair  $(V, E)$  is called a **graph**. We call elements of  $V$  the **vertices** and elements of  $E$  the **edges**. Sometimes we will write  $G = (V, E)$  and call the graph  $G$ .

When thinking about vertices, think about locations like the different locations (land masses and the island) in Königsberg, and when thinking about edges remember the bridges, each of which connected some pair of locations in the city. Alternatively, think of vertices as the train stations and think of each edge as the track between them. Notice that our abstract definition of a graph does not overtly have a visual component. However, we could draw a picture that corresponds to a graph by drawing a dot for each vertex and drawing a possibly curved line segment connecting a pair of vertices for each edge of the graph. As before, we would have to make certain that our representation clearly showed that edges do not intersect one another. Each edge is separate.

The word “graph” comes from the Greek root word meaning “to write”. In high school math classes, “to graph” means “to draw”, as in “graphing a

function”. Perhaps we should have chosen a different name since a graph is not inherently visual, but the term is too firmly entrenched to change now, and often an appropriate visual representation of a graph gives us valuable insights.

We can be somewhat satisfied with our definition, but now we have to step back and ask ourselves whether there are any issues that need to be addressed. If we look at the graph corresponding to the Königsberg Bridge Problem, we might notice a potential issue, namely, there are two land masses that are connected with two different bridges. In terms of the abstract definition of a graph, that means that the same pair of vertices appears as distinct edges in  $E$ . We have isolated an issue. So let’s explicitly allow  $E$  to contain multiple copies of a pair  $\{v, v'\}$ , just as we allowed multiple bridges between the same land masses in the Königsberg Bridge Problem. So we will allow multiple edges between the same pair of vertices and indicate their presence by writing the same pair down as many times as there are multiple edges between those vertices. After we have isolated the idea of multiple edges, we can define graphs with that feature.

*Definition.* A graph  $G = (V, E)$  is said to have **multiple edges** if  $E$  contains two (or more) distinct copies of an edge  $\{v, v'\}$ . In plain language,  $G$  has multiple edges if it has two vertices that are connected by more than one edge. Technically, the existence of multiple edges connecting the same two vertices means that  $E$  is a multiset, not a set, but we will ignore this issue.

We’ve gotten pretty abstract, pretty quickly. The following exercise is to make sure you’re following.

*Exercise 2.5.* Abstract the situation of a group of people shaking hands, each shaking hands with some or all of the other people, to a graph, where vertices correspond to people and edges to handshakes. What would it mean for this graph to have multiple edges?

Another issue that comes to mind is whether the edge  $\{v, v'\}$  is the same or different from the edge  $\{v', v\}$ . That is, does the order of the vertices in an edge make a difference? Well, we could choose either answer. In the situations that generated our concept, the order did not matter (we could walk over the bridges in either direction and the handshakes did not have a direction to them, for example), so we will choose not to distinguish between  $\{v, v'\}$  and  $\{v', v\}$ . So for this concept of a graph, we could replace any of our edges with a pair of vertices in the opposite order and say that that is the same graph.

If we chose to view differently ordered edges as different, then we would be describing something that is referred to as a *directed graph*. Directed

graphs would be appropriate for capturing some other situations. For example, suppose there were one-way signs on the bridges of Königsberg, then a directed graph would be required to capture the restrictions that the new problem presented. Directed graphs also make more sense when modeling the spread of a disease, since we would want the representation to capture the idea that an infected person infects a non-infected person.

We have yet one more issue that we may want to make a decision about: should we allow an edge to go from a vertex to itself? None of our generating scenarios has such a situation; however, we could easily imagine such a situation. We could imagine a bridge that starts and ends on the same land mass, like an overpass, for example. So we will choose to allow edges of the form  $\{v, v\}$ . Since that edge is rather distinctive looking, we will give it a name.

*Definition.* Let  $G = (V, E)$  be a graph with a vertex  $v$ . Then an edge of the form  $\{v, v\}$  is called a **loop** (at  $v$ ).

Now let's get accustomed to the vocabulary of a graph by looking at the Königsberg Bridge Problem in our new terms.

*Exercise 2.6.* Carefully, using the definitions we have just chosen, construct a graph for the Königsberg Bridge Problem,  $K = (V, E)$ . Give each vertex a name (probably just a letter); then, using these letters, write  $V$  and  $E$  for this graph.

*Exercise 2.7.* For each of the other challenges, think about how you would specify a graph  $G = (V, E)$ . It would be tedious to do all of them. So pick at least one to write out carefully.

In thinking about the Königsberg Bridge Problem, it would be reasonable to say that a bridge “has endpoints  $A$  and  $B$ ” where  $\{A, B\}$  was an edge in the graph. So there is some natural vocabulary that will help us to discuss questions about graphs.

*Definition.* Let  $G$  be a graph containing vertices  $v$  and  $v'$  and the edge  $e = \{v, v'\}$ . Then  $e$  has **endpoints**  $v$  and  $v'$ , and  $v$  and  $v'$  are **adjacent by**  $e$ .

Making this definition lets us use some more intuitive and familiar language to talk about graphs. In particular, now a graph has multiple edges if there it has a pair of vertices that are the endpoints of two distinct edges. However, there are some weird side effects too: if  $G$  contains the loop  $\{v, v\}$ , then  $v$  is adjacent to itself. If we're going to go to all the trouble to carefully create definitions, then we must also be careful when using common language to talk about the ideas.

When we look at our visual representations of the Königsberg Bridge Problem, the Paperperson's Puzzle, and the Gas-Water-Electricity Conundrum, one feature that we see in describing those graphs concerns the number of edges that emerge from each vertex.

*Definition.* If  $v$  is a vertex, then we define the **degree** of  $v$ , written as  $\deg(v)$ , to be the number of edges with an endpoint  $v$ , where a loop counts twice. The **total degree** of a graph  $G$  is the sum of the degrees of the vertices of  $G$ .

- Exercise 2.8.* 1. If a set  $E = \{\{v_1, v'_1\}, \{v_2, v'_2\}, \{v_3, v'_3\}, \dots, \{v_n, v'_n\}\}$  is the edges of a graph, how can we determine the degrees of the vertices without drawing the graph? (This notation just means that there are  $n$  edges in the graph; however, it does not tell us how many vertices there are. Any of these vertices could be the same as any other. For example,  $v_3$  might be the same as  $v_4$ . Or  $v_5$  could be the same as  $v'_5$ . But if you were given a specific set  $E$ , you would know which vertices were the same and which were different.)
2. Compute the degrees of the vertices in the Königsberg Bridge Problem using the procedure you described in the previous part of this exercise, and make sure those answers agree with the numbers you got by just looking at your visual representation.
3. Write out a specific example of a graph with at least five vertices and compute the degree of each vertex and the total degree of the graph.

It is now time for our first theorem. It points out that the total degree of any graph must be an even number.

*Theorem 2.9.* The total degree of a graph is even.

*Corollary 2.10.* Let  $G$  be a graph. Then the number of vertices in  $G$  with odd degree is even.

One of the habits of a good mathematician is to check how theorems work in particular cases every time you do a proof. This habit helps to make abstract mathematics meaningful.

*Exercise 2.11.* Confirm the truth of the theorem and corollary above in the Königsberg Bridge Problem graph and in the graph you constructed in part 3 of exercise 2.8.

Theorem 2.9 and its corollary point out restrictions on graphs with respect to the degrees of vertices. These insights allow us to determine whether graphs could exist with various properties.

*Exercise 2.12.* Determine whether the following data could represent a graph. For each data set that *can* represent a graph, determine *all* the possible graphs that it could be and describe each graph using pictures and set notation. If no graph can exist with the given properties, state why not.

1.  $V = \{v, w, x, y\}$  with  $\deg(v) = 2$ ,  $\deg(w) = 1$ ,  $\deg(x) = 5$ ,  $\deg(y) = 0$
2.  $V = \{a, b, c, d\}$  with  $\deg(a) = 1$ ,  $\deg(b) = 4$ ,  $\deg(c) = 2$ ,  $\deg(d) = 2$
3.  $V = \{v_1, v_2, v_3, v_4\}$  with  $\deg(v_1) = 1$ ,  $\deg(v_2) = 3$ ,  $\deg(v_3) = 2$ ,  $\deg(v_4) = 5$

You proved earlier that the total degree of any graph is even. Let's consider a sort of converse question, namely, if the orders of vertices are given such that the total degree *is* even, can we create a corresponding graph?

*Exercise 2.13.* If you are given a finite set  $V$  and non-negative integer for each element in the set such that the sum of these integers is even, can  $V$  be realized as the vertices of a graph with the associated degrees? If so, prove it. If not, give a counter-example.

The situation described in the Königsberg Bridge Problem was well modeled by the concept of a graph, which you have drawn. After abstracting the set-up, we must also translate the challenge of the problem in terms of the associated graph. In posing the Königsberg Bridge Problem, Otto was asking whether it is possible to trace every edge (bridge) of the graph without picking up the pencil and without going over any edge more than once.

*Exercise 2.14.* Try to trace your Königsberg Bridge graph without picking up your pencil and without going over any edge more than once. You can put the Sternbucks anywhere you like; try several locations. Does the starting place affect the answer?

If we can trace one visual representation of the Königsberg Bridge graph, we can trace any correct representation, which is why we can abuse language and talk about *the* (visual representation of the) graph when working on this problem. But to avoid this subtlety entirely, we can ask Otto's question about our graph, where the graph is presented in the set notation  $K = (V, E)$ .

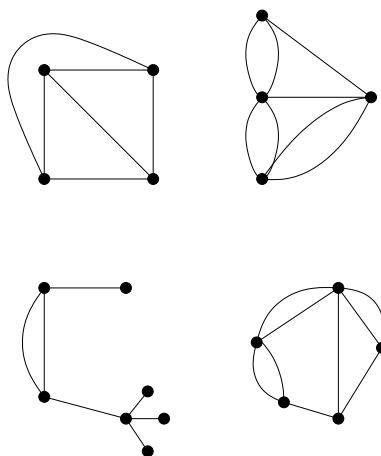
*Exercise 2.15.* Translate the Königsberg Bridge Problem into a question about its graph,  $K = (V, E)$ , without reference to a visual representation of  $K$ .

The Königsberg Bridge Problem was modeled by a graph and its challenge was described in terms of a tracing problem. This problem naturally encourages us to explore the general question of when we can trace a graph without picking up the pencil and without going over any edge more than once. To gain experience with this question, an excellent strategy is to try several graphs and observe which ones seem to be traceable as prescribed and which ones seem to be impossible. Then we can try to isolate what features of a graph seem to make it traceable.

*Exercise 2.16.* Draw the graph associated with the Paperperson's Puzzle. Try to trace the graph without picking up your pencil and without going over any edge more than once.

The more experience we get, the more apt we are to identify characteristics of a graph that indicate traceability.

*Exercise 2.17.* For each graph pictured below, try to trace the graph without picking up your pencil and without going over any edge more than once. Look for some feature or features among the graphs that you can trace compared to the ones that you can't. You may not be able to characterize those graphs that are traceable, but perhaps you can isolate some features of a graph that definitely make it traceable or definitely make it untraceable.



## 2.3 Taking a Walk

Looking at examples is a great way to begin to explore an idea, but at some point it is valuable to become a bit more systematic in the investigation.



Starting with simple cases is an excellent strategy for developing insight. So let's consider some simple graphs to see whether we can discover some sort of pattern among those that are traceable or untraceable.

Let's start with just one bridge (or edge). If there were only one bridge between two land masses, then the edge could be traced, but it would be impossible to return to the starting place without retracing the same edge. Recall that in the Königsberg Bridge Problem, Otto challenged Friedrich to return to his starting place, so we must consider that restriction. If only one bridge existed and it connected the same island to itself, then we could traverse the bridge while starting and ending at the same point. That is, if a graph had only one edge and that edge were a loop, then we could trace the graph.

Now let's consider graphs with two or three edges.

*Exercise 2.18.* Draw all possible graphs that contain two or three edges. Which are traceable and which are not traceable?

Now investigate graphs with four edges.

*Exercise 2.19.* Draw all graphs with four edges without loops. Which graphs with four edges are traceable? Try to be systematic and try to isolate some principles that seem pertinent to traceability.

Perhaps you will observe that the degrees of the vertices are important for the issue of traceability.

*Exercise 2.20.* For each of the graphs you drew in exercises 2.18 and 2.19 as well as those for the Königsberg Bridge Problem and the Paperperson's Puzzle, make a chart that records the degrees of each vertex of each graph. Do you see something that separates the good from the bad (traceable from not)?

We translated Otto's Königsberg Bridge Problem into a question about graph theory, and now we will formalize what it means to find a solution. Just as making our definitions for the abstraction process helped us decide what was important about the problem, formalizing the question helps us see how to break it down into more manageable steps.

The act of tracing the edges of a graph is a fairly clear process, but there are really several different ways of moving about a graph, some involving the proviso of not repeating edges and the more basic idea of just moving around. So let's take the step of pinning down some definitions about how we can move about on a graph. The first definition refers to moving from one vertex to another, but there is no restriction about repeating the same edge.

*Definition.* Let  $G$  be a graph with vertices  $v$  and  $w$ . A **walk** from  $v$  to  $w$ ,  $W$ , is a finite sequence of adjacent vertices and edges of  $G$  of the form

$$W: v(= v_0), e_1, v_1, e_2, v_2, e_3, \dots, v_{k-1}, e_k, w = (v_k)$$

where the  $v_i$ 's are vertices of  $G$ , and for each  $i$ ,  $e_i$  is the edge  $\{v_{i-1}, v_i\}$ . We explicitly allow a trivial walk from  $v$  to  $v$ ,  $T : v$ , which is just one vertex without any edges.

Since walks that do not repeat an edge are of special interest, for example in the Königsberg Bridge Problem, we give them a special name.

*Definition.* Let  $G$  be a graph with vertices  $v$  and  $w$ . A **path** from  $v$  to  $w$  is a walk,

$$P: v, e_1, v_1, e_2, v_2, e_3, \dots, v_{k-1}, e_k, w$$

where the edges  $e_i$  are all distinct. In other words,  $P$  is a walk that contains no repeated edges.

Another kind of walk we might consider is one in which we never visit the same vertex more than once.

*Definition.* A **simple path** from  $v$  to  $w$  is a path from  $v$  to  $w$ ,  $P : v(= v_0), e_1, \dots, e_n, w(= v_n)$ , such that  $v_i \neq v_j$  for  $i \neq j$ . In other words,  $P$  contains no repeated vertices.

The next theorem states that we really only need to check that a walk has no repeated vertices in order to conclude that it is a simple path, that is, it is unnecessary to first check that it is a path.

*Theorem 2.21.* Let  $G$  be a graph and  $W: v_0, e_1, \dots, e_n, v_n$  a walk in  $G$  such that the vertices  $v_i$  are all distinct. Then  $W$  is a simple path.

In the statement of the Königsberg Bridge Problem, recall that Otto's challenge involved returning to the starting place. So, just as we had a hierarchy of definitions associated with ways to get from one vertex to another, we have a similar hierarchy of definitions of types of perambulations that return us to the vertex from which we started. Notice that in our original definition of a walk, the beginning and ending vertices had no restrictions, so they could actually have been the same vertex.

*Definition.* Let  $G$  be a graph.

1. A **closed walk** is a walk  $W: v_0, e_1, \dots, e_n, v_n$  where  $v_0 = v_n$ . In other words,  $W$  begins and ends at the same vertex.
2. A **circuit** is a closed walk (with at least one edge) that does not contain a repeated edge.

3. A **simple circuit** is a circuit that does not have any repeated vertices except the initial/terminal vertex.

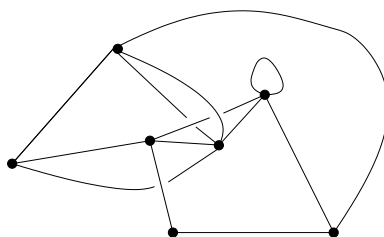
*Exercise 2.22.* Write out some of the circuits in  $K = (V, E)$ , the Königsberg Bridge Problem graph. Find at least one circuit that is *not* a simple circuit.

When we have described a mathematical entity, such as a graph, it is often useful to look at smaller such objects that are contained in it. This strategy leads to the concept of a subgraph.

*Definition.* Let  $G = (V, E)$  and  $G' = (V', E')$  be graphs. If  $V' \subset V$  and  $E' \subset E$ , then we say that  $G'$  is a **subgraph** of  $G$ .

Much like checking a theorem in special cases to understand its meaning more thoroughly, mathematicians observe definitions in examples to help them understand the definitions.

*Exercise 2.23.* Consider the graph  $G$  below. Draw all subgraphs of  $G$  with three vertices.



Note that a circuit in a graph is *not* a subgraph; however, the set of vertices and the set of edges in the circuit do form a subgraph.

When talking about a vertex in a graph and also in a subgraph, confusion could arise about the vertex's degree. If  $G'$  is a subgraph of  $G$ , and a vertex,  $v$ , is in both  $G$  and  $G'$ , then we will use the notation  $\deg_G(v)$  and  $\deg_{G'}(v)$  to denote its degree in  $G$  and  $G'$  respectively.

*Theorem 2.24.* Let  $G$  be a graph. Let  $C$  be a subgraph of  $G$  that consists of the vertices and edges that belong to a circuit in  $G$ . Then  $\deg_C(v)$  is even for every vertex,  $v$ , of  $C$ .

One basic question we can ask about a graph is whether we can get from one vertex to another.

*Definition.* Let  $G$  be a graph with vertices  $v$  and  $w$ .

1. We say  $v$  is **connected to**  $w$  if there exists a walk from  $v$  to  $w$ .

2. The graph  $G$  is **connected** if every pair of vertices of  $G$  is connected.  
If not, we say  $G$  is **disconnected**.

*Exercise 2.25.* Show that your graph of the Königsberg Bridge Problem is connected. Carefully use the definitions. Also, give an example of a graph that is not connected.

It is obvious, visually, when a graph is connected, at least when it has a small number of vertices, but that is different from a proof. As the last exercise hopefully showed you, there's a lot to write down to show that a graph is connected. The following theorem helps shorten the work; it also tells us that the term "connected" behaves as we use it in common English.

*Theorem 2.26.* Let  $G$  be a graph with vertices  $u$ ,  $v$ , and  $w$ .

1. The vertex  $v$  is connected to itself.
2. If  $u$  is connected to  $v$  and  $v$  is connected to  $w$ , then  $u$  is connected to  $w$ .
3. If  $v$  is connected to  $w$ , then  $w$  is connected to  $v$ .

This lemma lets us define the *connected components* of  $G$  as the subgraphs of mutually connected vertices with all their edges.

*Definition.* Let  $G = (V, E)$  be a graph. Then a subgraph of  $G$ ,  $H = (V', E')$ , is a **connected component of  $G$**  if every vertex that is connected to a vertex in  $V'$  by a walk in  $G$  is already in  $V'$  and every edge from  $E$  with endpoints in  $V'$  is already in  $E'$ .

*Exercise 2.27.* Construct a graph that has more than one connected component. What can you say about the traceability of this graph?

The Königsberg Bridge Problem produced a graph that we sought to traverse without lifting our pencil or repeating an edge. So that problem gives rise to a definition that captures this kind of traceability.

*Definition.* Let  $G$  be a graph with vertices  $v$  and  $w$ .

1. An **Euler circuit** for  $G$  is a circuit in  $G$  that contains every vertex and every edge of  $G$ .
2. An **Euler path** from  $v$  to  $w$  is a walk from  $v$  to  $w$  that contains every vertex and contains every edge exactly once.

All of these definitions concerning ways to get around on a graph were motivated by trying to capture ideas suggested in the Königsberg Bridge Problem. So let's see whether we can restate that puzzle using our new vocabulary.

*Exercise 2.28.* Restate the Königsberg Bridge Problem using our formal definitions.

In some sense, restating a question in formal terms does not make any progress towards solving it; however such a restatement can be helpful. Now we are clear on what we seek to find in our graph: we seek a walk with various restrictions. Getting in the habit of using the definition of a walk (that is, a sequence of vertices and edges) can be helpful in proving theorems about graphs such as the next one, which shows that if we can get from one vertex to another in a graph, we can go between the two without repeating edges or vertices. Remember to use the definitions of each term. If you are uncertain what level of detail is required for a rigorous proof, imagine that your reader is a computer that will use your writing as instructions. In the following theorem, the computer would start with the data of a connected graph (what does that mean?), and you should tell it how to build a simple path between an arbitrary pair of vertices.

*Theorem 2.29.* If  $G$  is a connected graph, then any two vertices in  $G$  can be connected by a simple path.

One natural attempt to solve the Königsberg Bridge Problem would be to simply start walking without going over the same bridge twice and see where you end up. This approach somewhat reflects our strategy of formalizing Otto's challenge. We constructed the definition of an Euler circuit from ideas of walks and paths, so we might consider building an Euler circuit from those same concepts.

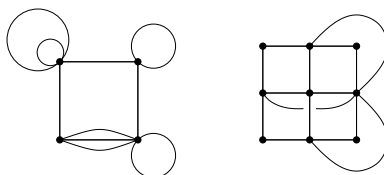
We are now ready to characterize those graphs that have an Euler circuit or Euler path. Determining whether a graph has an Euler circuit or Euler path turns out to be easy to check. It is always a pleasure when a property that appears to be difficult to determine actually is rather simple.

*Theorem 2.30 (Euler Circuit).* A graph  $G$  has an Euler circuit if and only if it is connected and every vertex in  $G$  has even degree.

If you truly understand the proof of this theorem, you should be able to take a graph and produce an Euler circuit, if it has one, using the technique implicit in the proof of this theorem. So here is an exercise that lets you explore the method of the proof rather than just the statement of it.

*Exercise 2.31.* In the following graphs, find an Euler circuit using the method that successfully proved the Euler Circuit Theorem.

We can now definitively complete the Königsberg Bridge Problem by translating our solution back into the language of bridges and locations.



*Exercise 2.32.* Solve the Königsberg Bridge Problem. Write your solution in a way that Otto could understand from start to finish, that is, write your answer thoroughly in ordinary English, not Old Prussian.

Similarly, we can settle the Paperperson's Puzzle.

*Exercise 2.33.* Solve the Paperperson's Puzzle.

We've finished with Otto's challenge and the Königsberg Bridge Problem, but now we need to think about what other kinds of theorems are true. The first place to look for new theorems is in modifying theorems we've already proven. The second place to look is back at the actual proofs we've produced; sometimes when looking back and summarizing an old proof we realize that simply changing the hypotheses would produce new theorems, or that we've actually proven something more than we set out to show.

In that vein, we can ask: how much easier is it to trace a graph if we don't have to end where we started? It turns out to be only a little easier.

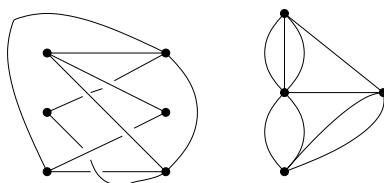
*Theorem 2.34 (Euler Path).* A graph  $G$  has an Euler path if and only if  $G$  is connected and has zero or two vertices of odd degree.

Let's make certain that the distinction between an Euler path and an Euler circuit is clear.

*Exercise 2.35.* Give an example of a graph with an Euler path but not an Euler circuit. What must be true of any such example?

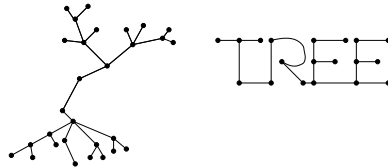
Again, let's practice the method of proof for the Euler Path Theorem.

*Exercise 2.36.* In the following graphs, find an Euler path using the method that successfully proved the Euler Path Theorem.



## 2.4 Trees

We've proven a large number of theorems about graphs with circuits and when graphs have certain kinds of circuits. We now turn our attention to some interesting theorems about graphs without circuits, *trees*.



*Definition.* 1. A graph is called a **tree** if it is connected and has no circuit.

2. A subgraph that is a tree will be called a **subtree**.

*Exercise 2.37.* A start-up airline, AirCheap, only flies to four cities, and all flights go through Wichita. But from Wichita you can fly to Austin, Denver, or Chicago. Construct a graph that has vertices corresponding to the cities and edges corresponding to flights for AirCheap. Is the graph a tree? Justify your answer.

One feature of a tree is that it must contain an edge, one end of which is unattached to other edges. These endpoints with degree one are sometimes called **leaves**.

*Theorem 2.38.* Any tree that has more than one vertex has a vertex of degree one, in fact, it has at least two vertices of degree one.

When we actually look at a tree, we notice that there are often quite a few vertices of degree one (leaves). This observation invites us to explore the question of how the number of degree one vertices relates to other features of the tree.

*Exercise 2.39.* By looking at a few examples, explore the relationship between the number of degree one vertices of a tree and other features of the tree. Make a conjecture and prove it.

We can tell whether a graph is a tree simply by comparing the number of its vertices with the number of its edges.

*Exercise 2.40.* There is a simple relationship between the number of vertices and edges in a tree. Make a conjecture of the following form and prove it: A graph with  $n$  vertices is a tree if and only if  $G$  is connected and has \_\_\_\_\_ edges.

Trees are particularly simple examples of graphs. In graphs with circuits, there are often many different ways to get from one vertex to another, but in a tree, there is only one option.

*Theorem 2.41.* If  $v$  and  $w$  are distinct vertices of a tree  $G$ , then there is a unique path in  $G$  from  $v$  to  $w$ .

This implies that trees are disconnected by the removal of any edge.

*Corollary 2.42.* Suppose  $G = (V, E)$  is a tree and  $e$  is an edge in  $E$ . Then the subgraph  $G' = (V, E \setminus \{e\})$  is not connected.

On the contrary, removing an edge from a circuit in a graph will not disconnect the graph. Of course, such a graph is definitely not a tree.

*Theorem 2.43.* Let  $G = (V, E)$  be a connected graph that contains a circuit  $C$ . If  $e$  is an edge in  $C$ , then the subgraph  $G' = (V, E \setminus \{e\})$  is still connected.

Every graph has subtrees, and connected graphs have subtrees that contain all the vertices of the graph. Sometimes we can use these subtrees as starting points for analyzing the larger graph.

*Theorem 2.44.* Let  $G$  be a connected graph. Then there is a subtree,  $T$ , of  $G$  that contains every vertex of  $G$ .

*Definition.* Let  $G$  be a graph. Then a subtree  $T$  of  $G$  is a **maximal tree** if and only if for any edge of  $G$  not in  $T$ , adding it to  $T$  produces a subgraph that is not a tree. More formally, a subtree  $T = (W, F)$  of  $G = (V, E)$  is a **maximal tree** if and only if for any  $e = \{v, w\}$  in  $E \setminus F$ ,  $T' = (W \cup \{v\} \cup \{w\}, F \cup \{e\})$  is not a subtree.

*Theorem 2.45.* A tree  $T$  in a connected graph  $G$  is a maximal tree if and only if  $T$  contains every vertex of  $G$ .

## 2.5 Planarity

Earlier, we ran across the issue of whether we could draw a graph in the plane without having edges cross. If a graph can be drawn without edges crossing, we can often use geometric insights to deduce features about the graph. In the next two sections, we investigate issues concerning graphs that can be drawn in the plane without edges crossing.

*Definition.* A graph  $G$  is called **planar** if it can be drawn in the plane ( $\mathbb{R}^2$ ) such that the edges only intersect at vertices of  $G$ .

Remember that a graph is just two sets  $G = (V, E)$ . When a graph is presented in this formal way, it is far from obvious whether the graph



is planar. To aid in our exploration of planarity, let's describe some new families of graphs.

- Definition.*
1. For a positive integer  $n$ , the **complete graph on  $n$  vertices**, written  $K_n$ , is the graph having  $n$  vertices and containing no loops such that each pair of vertices is connected by a unique edge.
  2. For positive integers  $m$  and  $n$ , the **complete bipartite graph**,  $K_{m,n}$ , is the graph having  $m+n$  vertices, each of the first  $m$  vertices connected to each of the last  $n$  vertices by a unique edge and having no other loops or edges.

*Exercise 2.46.* Draw graphs of  $K_3$ ,  $K_4$ ,  $K_5$ ,  $K_{2,3}$ ,  $K_{3,3}$ , and  $K_{2,4}$ . Which appear to be planar graphs? Are any of them familiar?

The next theorems are quite hard to prove rigorously. Showing that something is planar only requires finding one particular way to draw it; but showing that something is not planar involves showing that no arrangement is possible. Instead of trying to find ironclad proofs now, give informal, plausible arguments that they are true. Later we will be in a position to give firmer proofs.

*Theorem\** 2.47. The graph  $K_{3,3}$  is not planar.

*Theorem\** 2.48. The graph  $K_5$  is not planar.

Even if we can't prove these theorems, we can interpret their consequences. Recall the Gas-Water-Electricity Dilemma, namely: Three new houses have just been built in Houseville, and they all need natural gas, water, and electricity lines, each of which is supplied by a different company. Can the connections be made without any crossings?

*Exercise 2.49.* What do the previous theorems imply about the Gas-Water-Electricity Dilemma and the Five Station Quandry?

On the other hand, trees can always be drawn in the plane.

*Theorem* 2.50. Let  $G$  be a tree. Then  $G$  is planar.

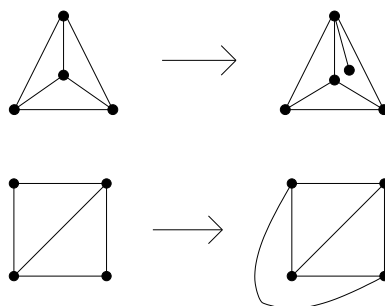
It is difficult to decide what makes a graph planar without considering non-planar graphs. We've run into two examples of non-planar graphs thus far:  $K_{3,3}$ , the graph describing the Gas-Water-Electric Dilemma, and  $K_5$ , the graph that represents the Five Station Quandry. If we take away any one edge from either of these graphs, we produce planar subgraphs.

*Exercise 2.51.* Show that if we remove any one edge from either  $K_{3,3}$  or  $K_5$ , the resulting subgraphs are planar.

## 2.6 Euler Characteristic

As we've mentioned before, sometimes hard facts can be proven by starting with simple cases and building up to more complex situations. Having control of the different ways that we can build more complex situations makes this technique even more powerful. It turns out that every connected, planar graph can be drawn from the simplest graph (one vertex with no edges) by repeatedly taking one of two steps.

*Theorem 2.52* (Constructing Planar Graphs). Let  $G$  be a connected, planar graph. Then  $G$  can be drawn in a finite number of steps using only the following two procedures. At each stage, the procedure will produce a connected subgraph with a planar drawing. Start with a subgraph,  $G_0$ , that has only one of the vertices from  $G$  and no edges, with the obvious planar drawing. Given a planar drawing of the subgraph  $G_n$ , build a planar drawing of the next one,  $G_{n+1}$ , by doing one of the following procedures. (1) For some vertex  $v$  in  $G$  not in  $G_n$  that is adjacent to a vertex  $v'$  in  $G_n$  by the edge  $e = \{v, v'\}$ , add  $v$  and the edge  $e$  to the drawing of  $G_n$ . (2) Draw an edge of  $G$  whose endpoints are two vertices that are already in  $G_n$ .

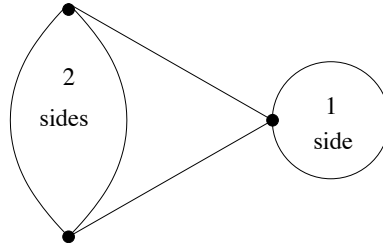


In fact, we can improve the above theorem by specifying the order in which we carry out the two procedures.

*Theorem 2.53.* Let  $G$  be a connected, planar graph. Then  $G$  can be drawn in the plane first doing all the procedure (1) steps and then doing all the procedure (2) steps.

A graph drawn in the plane chops  $\mathbb{R}^2$  into a number of regions. We will call these regions **faces**, and we will include the region “outside” the graph, called the unbounded region, as one of the faces. Each face is bounded by edges of the graph. Notice that the face inside a loop has only one side if no other part of the graph is drawn inside the loop. And the face between a pair of multiple edges has only two sides if no other part of the graph is

drawn between the two edges. Even weirder, the simplest graph, which has one vertex and no edges, has one face with *no* sides.



For any graph  $G = (V, E)$ , let  $|V|$  denote the number of vertices of  $G$  and  $|E|$  denote the number of edges of  $G$ . These numbers do not depend on a drawing of  $G$ . However, it turns out that every drawing of a planar graph in the plane will have the same number of faces as well. So if  $G$  is planar with a fixed drawing in the plane, let  $|F|$  denote the number of faces in that drawing of  $G$ . The fact that  $|F|$  does not depend on the drawing of  $G$  is quite surprising from the perspective of our definition of a graph, and we will prove it shortly. For now, let's just check this assertion with some examples.

*Exercise 2.54.* Draw a planar graph with at least five vertices and five faces. Now produce another planar drawing of the same graph that is as different as you can make it. Compare the number of faces in each drawing.

When we begin to draw a planar graph in the Constructing Planar Graphs Theorem, we start with a single vertex, no edges, and one face. As we add edges, using the two procedures in the Constructing Planar Graphs Theorem, we produce graphs that have different numbers of vertices, edges, and faces. By investigating how these two procedures change  $|V|$ ,  $|E|$ , and  $|F|$ , we are able to say something about how these numbers are related.

*Exercise 2.55.* Draw a graph using the two procedures detailed in the Constructing Planar Graphs Theorem. Create a chart that includes the number of vertices, number of edges, and number of faces at each stage. Do you notice any patterns?

If you were successful with the preceding exercise, you will have discovered one of the most famous formulas in graph theory.

*Theorem 2.56* (Euler Characteristic Theorem). For any connected graph  $G$  drawn in the plane,

$$|V| - |E| + |F| = 2.$$

The Euler Characteristic Theorem allows us to deduce the result about the invariance of the number of faces. Notice that this next corollary does not require that the planar graph  $G$  be connected.

*Corollary 2.57.* Let  $G$  be a planar graph. Then any two drawings of  $G$  in the plane have the same number of faces.

The Euler Characteristic Theorem also gives us a new proof of an old fact.

*Corollary 2.58.* A graph  $G$  with  $n$  vertices is a tree if and only if  $G$  is connected and has  $n - 1$  edges.

The Euler Characteristic Theorem has many consequences including some theorems about the relationship between the numbers of vertices and edges of a connected planar graph.

*Theorem 2.59.* Let  $G$  be a connected, planar graph with no loops or multiple edges having  $|V|$  vertices and  $|E|$  edges. If  $|V| \geq 3$ , then  $|E| \leq 3|V| - 6$ .

If the theorem seems elusive, try this lemma first.

*Lemma 2.60.* Let  $G$  be a planar graph with no loops or multiple edges. Then  $G$  can be realized as a subgraph of a planar graph  $H$ , drawn such that all faces of  $H$  have three sides.

In general, when we consider a graph it may be difficult to prove for certain that it is impossible to draw it in the plane. How do we know that we simply haven't thought of some clever way to draw it? Conditions like those in the previous theorem on the relationship between the numbers of vertices and edges in a connected planar graph can be used to show us that certain graphs are not planar.

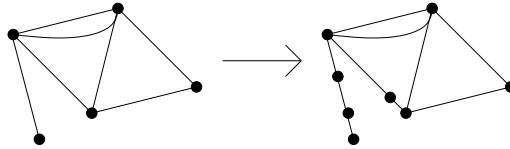
*Corollary 2.61.* The graph  $K_5$  is not planar.

A little more analysis is required to prove that  $K_{3,3}$  is not planar.

*Theorem 2.62.* The graph  $K_{3,3}$  is not planar.

Clearly, if we have a graph built from  $K_5$  or  $K_{3,3}$  by adding vertices and edges, it cannot become planar, because if we could draw the bigger graph in the plane, then that would put  $K_5$  or  $K_{3,3}$  in the plane. Also, adding extra vertices in the middle of edges does not affect the planarity of a graph. This observation leads to the following definition.

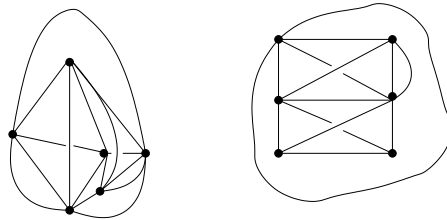
*Definition.* A graph  $G' = (V', E')$  is a **subdivision of a graph**  $G$  if  $G'$  is obtained from  $G = (V, E)$  by adding a new vertex  $u$  to  $V$  and replacing an edge  $\{v, w\}$  with two edges  $\{v, u\}$  and  $\{u, w\}$  and repeating this process a finite number of times. Graphically, a subdivision  $G'$  of  $G$  is simply built by inserting zero or more vertices into the interiors of edges of  $G$ .



The following theorem completely characterizes whether a graph is planar or not. It turns out that planarity of graphs hinges entirely on the specific graphs  $K_{3,3}$  and  $K_5$ , the graphs we know as the Gas-Water-Electricity graph and the Five Station Quandry graph. Unfortunately, the following theorem is difficult to prove.

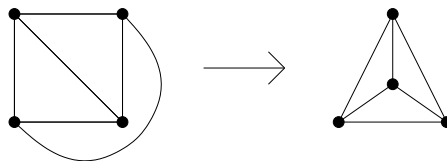
*Theorem\* 2.63 (Kuratowski).* A graph  $G$  is planar if and only if  $G$  contains no subgraph that is a subdivision of  $K_{3,3}$  or  $K_5$ .

*Exercise 2.64.* For each of the following graphs, find a subgraph that is a subdivision of  $K_{3,3}$  or  $K_5$  or find a way to draw it that demonstrates that it is planar.



One interesting feature of graphs is that if they can be drawn in the plane at all, they can be drawn there with straight edges.

*Theorem\* 2.65.* Let  $G$  be a planar graph with no loops or multiple edges. Then  $G$  can be drawn in the plane in such a way that every edge is straight.



## 2.7 Colorability

As we said when we started, one of the reasons to abstract a problem is that the techniques and concepts that we create when solving one problem may

help us in other situations. Graph theory captures connectivity and adjacency, so questions that use these terms might benefit from graph theoretic insights.

For example, have you ever wondered how “they” pick the colors for the countries or states on a map or globe? Well, one requirement is that adjacent countries have different colors. Under the constraint that adjacent countries have different colors, how many colors are necessary to color a map? It is the use of the word “adjacent” here that makes us think that graph theory might be useful in attacking this question.

First we need to abstract the problem and find a graph somewhere. Let’s work with the continental United States for now. There are at least two natural ways to make a map into a graph. The first is to just make the state borders into the edges and the intersections of multiple state borders (like the Four Corners point) into the vertices. Then the problem has something to do with coloring the bounded regions, the faces. We’ve mostly been talking about edges and vertices thus far, so this formulation is vaguely unsatisfying.

On top of that, when describing the problem, we said that adjacent states needed to be different colors. The only time the word adjacent appears in our graph theory definitions is for the endpoints of an edge. Recall (or look up) how we defined “adjacent”. To use this language, we need the states to be vertices in a graph. Two states are adjacent if they share a border, so we need edges between bordering states (vertices).

Notice that we could do either of these procedures to draw a graph right on the map without having edges intersect, so in either case the resulting graphs are planar. At first glance, these graphs look quite different. But they both contain the information about which states are next to which others. So there should be a way to build one graph from the other without thinking about the original map.

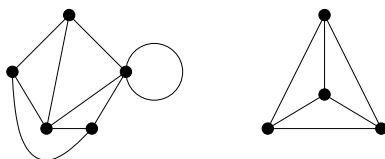
*Exercise 2.66.* Find a map of (a portion of) the United States that you can draw on, something pretty large. Construct graphs from this map using the two different procedures detailed above. Use different colored pens for the two graphs so that the two graphs are clearly visible. Describe how the two graphs are related.

These two graphs in the plane are called **dual graphs**, or more precisely, each is the dual graph of the other.

*Definition.* Let  $G = (V_G, E_G)$  be a planar graph with a fixed planar drawing. Construct a new graph  $\hat{G} = (V_{\hat{G}}, E_{\hat{G}})$  as follows. For each face in the drawing of  $G$ , including the unbounded one, add a vertex to  $V_{\hat{G}}$ . Notice that each

edge in  $G$  is a side of two faces in the drawing of  $G$ . So for each edge  $e$  in  $G$ , add an edge to  $E_{\hat{G}}$  connecting the two faces for which  $e$  was a side in the planar drawing of  $G$ . We will call  $\hat{G}$  the **dual of  $G$** . Note that this process can be done in the same plane as the drawing of  $G$  without crossing edges, so  $\hat{G}$  will also be planar (with a particular planar drawing).

*Exercise 2.67.* Draw dual graphs for each of the following graphs. Then construct the dual graph for each of the graphs you've constructed. Notice anything interesting at either step?



The fact that our two procedures for drawing on maps produce dual graphs means that each one can be used to produce the other. So any information that we can glean from one can be gleaned from the other. In other words, studying one is fundamentally the same as studying the other. In particular,  $|V_G| = |F_{\hat{G}}|$ ,  $|F_G| = |V_{\hat{G}}|$ , and  $|E_G| = |E_{\hat{G}}|$ .

Well, hopefully you're convinced that either procedure will do as a starting point for our abstraction from a map and that any fact we can prove about one tells a corresponding fact about its dual. The problem of coloring maps is usually considered by taking a vertex for each state and an edge between any two vertices in states that share a border. If we return to the map coloring problem, we must translate our challenge to refer to this new graph that we have created. In this representation, the problem asks us to assign a color to each vertex such that adjacent vertices are different colors.

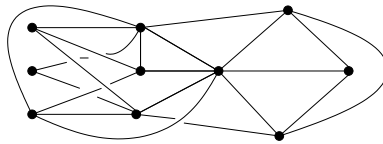
*Definition.* An  $n$ -**coloring** of a graph is a fixed assignment of a color to each vertex such that adjacent vertices are not the same color using at most  $n$  colors. A graph is  $n$ -**colorable** if it has an  $n$ -coloring.

Notice that the definition of  $n$ -coloring can refer to graphs that are not planar.

*Exercise 2.68.* Is the following graph 6-colorable? What is the smallest  $n$  such that this graph is  $n$ -colorable?

By the way, we should note that states that just touch at one point do not share a border. Any number of states could come together at one point.

Notice that the map coloring problem is simple for maps with a small number of states. Certainly, if we wanted to color a map with only 5 colors



and if there were only 5 or fewer states on our map, it would be easy. Just color each state a different color and then no state shares a border with a similarly colored state. The remainder of the section is dedicated to proving that 5 colors are enough to color any planar graph, that is, 5 colors are sufficient to color the vertices of any planar graph such that no two adjacent vertices have the same color. Our strategy for proving this theorem is to isolate conditions under which it is possible to extend a 5-coloring of a subgraph to a 5-coloring of a larger graph.

*Theorem 2.69.* Consider a graph,  $G$ , that is built from a subgraph,  $H$ , by adding one new vertex,  $v$ , and new edges that connect the new vertex to vertices in  $H$ . If we can color the subgraph  $H$  with five colors such that the new vertex,  $v$ , is not adjacent to vertices of all five colors, then we can color  $G$  with five colors.

One circumstance under which a new vertex will not be adjacent to vertices of all five colors is when the new vertex is not adjacent to five vertices altogether.

*Lemma 2.70.* If a graph  $G$  has no loops and is the union of a 5-colorable subgraph,  $H$ , and a new vertex,  $v$ , with edges such that  $v$  has  $\deg_G(v) < 5$ , then  $G$  is 5-colorable.

The combination of this last theorem and lemma suggest an inductive approach to answering the map coloring problem. But sadly, not all planar graphs have a vertex of degree less than 5.

*Exercise 2.71.* Construct a planar graph with no loops or multiple edges that contains no vertex of degree less than 5.

Although not all planar graphs have vertices of degree less than 5, they do have vertices of degree less than or equal to 5.

*Theorem 2.72.* Any planar graph with no loops or multiple edges has a vertex of degree at most 5.

Hint, look back at theorems about planar graphs.

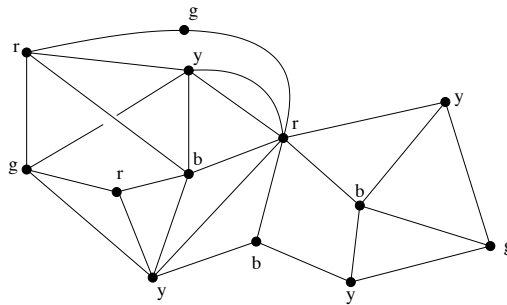
*Exercise 2.73.* Theorem 2.72 requires all of its hypotheses, of which there are three. For each hypothesis, find a counter-example to the theorem if that hypothesis were removed.



A graph might have several different  $n$ -colorings, and selecting a good one could be useful. Another way to think about this is to ask how we can change a fixed coloring to produce a new one with desired properties.

Let  $H$  be a 5-colorable graph with a fixed coloring, and let  $S$  be a subset of the colors. Then define  $H_S$  to be the subgraph of  $H$  that contains all of the vertices with colors in  $S$  and all of the edges both of whose vertices are vertices with colors in  $S$ .

*Exercise 2.74.* For the graph  $H$  below, which has been colored with the colors  $\{r, b, y, g\}$ , construct  $H_{\{r,y\}}$ ,  $H_{\{r,b,g\}}$ , and  $H_{\{y\}}$ .



*Lemma 2.75.* Let  $G$  be a graph without loops that is the union of a 5-colorable subgraph,  $H$ , and a new vertex  $v$  of degree 5. Fix a 5-coloring of  $H$  and label the five adjacent vertices to  $v$  as  $v_1, v_2, v_3, v_4, v_5$  with colors  $c_1, c_2, c_3, c_4, c_5$  respectively. Suppose that  $v_i$  and  $v_j$  are not connected in  $H_{\{c_i, c_j\}}$  for some pair of vertices/colors. Then  $G$  is 5-colorable.

Proving the preceding lemma involves writing down a procedure for finding a 5-coloring of  $G$  given the hypotheses. Is your procedure written so that a person or computer could use it to actually find the 5-coloring of  $G$ ? If not, you're not done yet.

Using the fact that  $G$  is planar, we will show that at least one of two special two-color subgraphs must be disconnected.

*Lemma 2.76.* Let  $G$  be a planar graph without loops or multiple edges that is the union of a 5-colorable graph  $H$  and a new vertex  $v$  of degree 5. Fix a 5-coloring of  $H$  and label the five adjacent vertices to  $v$  in cyclic order around  $v$  as  $v_1, v_2, v_3, v_4, v_5$  with colors  $c_1, c_2, c_3, c_4, c_5$  respectively. Then either  $v_1$  and  $v_3$  are not connected in  $H_{\{c_1, c_3\}}$ , or  $v_2$  and  $v_4$  are not connected in  $H_{\{c_2, c_4\}}$ .

Use the lemmas above to prove the following 5-color theorem.

*Theorem 2.77* (Five Color Theorem). Any planar graph with no loops is 5-colorable.

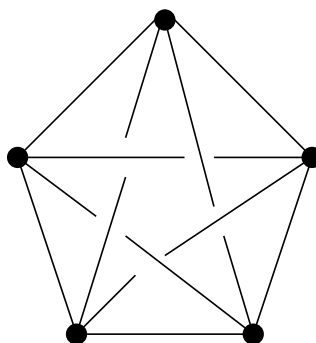
*Corollary 2.78.* Any map with connected countries can be colored with five colors such that no two countries that share a border have the same color.

In fact, four colors suffice to color any map. The following Four Color Theorem was a famous unsolved problem for more than a hundred years before it was proven using exhaustive computer methods. Its proof uses the Euler Characteristic Theorem extensively as well as techniques like those that we developed for switching colorings in  $H$  above, but the proof is extremely complicated and feels a bit unsatisfying in that the proof involves many cases that can be checked only by computers.

*Theorem 2.79* (Four Color Theorem). Any planar graph with no loops is 4-colorable.

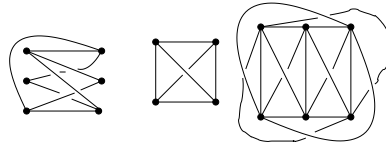
## 2.8 Regular Planar Graphs

In some sense, the Planarity section was about finding a better picture of a graph. If we can draw it in the plane, then it really lives in a simpler place than a non-planar graph. In a similar vein, mathematicians like pictures that capture symmetry. For example, there should be ways of drawing  $K_5$  that show that each of the 5 vertices is exactly the same as the others—there is no difference in the vertices in terms of the graph data.



A connected, planar graph  $G$  is said to have a **symmetric planar drawing** if it has a planar drawing where all of its vertices have the same degree and each face is bounded by the same number of edges.

*Exercise 2.80.* Consider the three graphs below and determine which ones have a symmetric planar drawing and which do not. Justify your answers.



We say a graph is a **regular planar graph** if it has a symmetric planar drawing where each vertex has degree at least 3, each face has at least 3 sides, and each edge bounds two distinct faces. Notice, for example, that the center graph in the previous problem is an example of a regular planar graph. Mathematicians love it when we require an object to have some property like symmetry and we're led to a finite list of possibilities.

*Exercise 2.81.* Find 5 regular planar graphs and prove that your collection is complete. (Hint: Let's denote the number of sides forming the boundary of each region in the plane by  $s$  and the degree of the vertices by  $d$ . Now express  $|F|$ , the number of faces, in terms of  $|E|$ , the number of edges, and  $s$ . Also express  $|V|$ , the number of vertices, in terms of  $|E|$  and the vertex degree,  $d$ .)

The previous exercise allows us to prove one of the central facts about symmetrical solids called the regular solids. A regular solid (also called a Platonic solid) is a convex, solid object with flat faces such that every face has the same number of edges and every vertex has the same degree.

*Theorem\** 2.82 (Platonic Solids Theorem). There are only five regular solids.

## 2.9 Morphisms

One of the interesting things about regular planar graphs is their high degree of symmetry: there are many ways to permute the names of the vertices that produce exactly the same graph. For a permutation of the vertices of a graph to have a chance at giving you back your original graph, each vertex must be moved to another vertex with the same degree. That condition is not sufficient to make a permutation of the vertices correspond to a symmetry of the graph. Let's pin down what we should mean by a symmetry of a graph.

*Definition.* 1. Let  $G = (V, E)$  and  $G' = (V', E')$  be graphs and let  $\phi : V \rightarrow V'$  be a function. Then we can extend  $\phi$  to  $E$  in the following way. If  $\phi(v) = v'$  and  $\phi(w) = w'$ , then we can set  $\phi(\{v, w\}) = \{v', w'\}$ . This new object may or may not be an edge of  $G'$ , but if it is for every edge in  $E$ , then we say that  $\phi$  is a **morphism** of graphs.

2. If  $\phi : G \rightarrow G'$  is a morphism of graphs that is a bijection of  $V$  with  $V'$  and  $E$  with  $E'$ , then we say  $\phi$  is an **isomorphism** of graphs.
3. If  $\phi$  is an isomorphism of graphs from  $G$  to itself, we call it an **automorphism** of graphs.

These definitions allows us to be specific about what we mean when we say that two graphs are the same. A graph  $G = (V, E)$  should be the same as  $G' = (V', E')$  if  $V'$  is just a relabeling of  $V$  and  $E'$  is the corresponding relabeling of  $E$ . This correspondence is exactly what the definition of an isomorphism of graphs captures. For this chapter, your intuitive understanding of when two graphs are the same has certainly been sufficient. A graph automorphism captures the idea of a symmetry of the graph, since the automorphism takes the graph to itself in a one-to-one manner.

Automorphisms are a great transition to groups, which is our next topic. The next exercise helps you to process these definitions and hopefully will provide a source of rich examples when the ideas return.

*Exercise 2.83.* How many automorphisms does  $K_5$  have? How many automorphisms does  $K_{3,3}$  have?

*Exercise 2.84.* For each of the 5 regular planar graphs, find all automorphisms of the graph.

## Chapter 3

# Group Theory

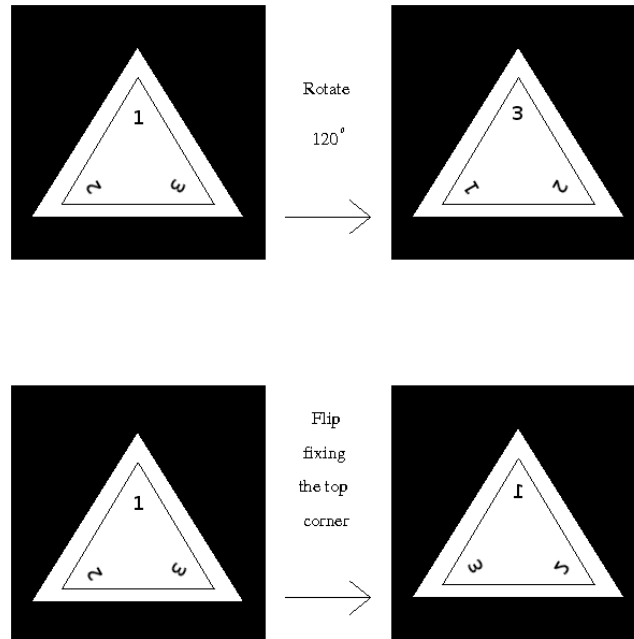
### 3.1 Examples Lead to Concepts

One of the most powerful and effective methods for creating new ideas is to look at familiar parts of our world and isolate essential ingredients. In mathematics, this strategy is particularly effective when we find several familiar examples that seem to share common features. So we will begin our next exploration by looking carefully at *adding*, at *multiplying*, and at *moving blocks*, with an eye toward finding similarities.

(1) Adding: Among the first computational skills we learn in our youths is addition of integers. So our first example is the familiar integers accompanied, as they are, by the method of combining them through addition.

(2) Multiplying: Real numbers are among our next mathematical objects and multiplication is another method of combining a pair of numbers to produce another number.

(3) Moving blocks: This example involves an equilateral triangular block fitting into a triangular hole, presenting challenges that you might recall from the first years of your life. As an inquisitive toddler, you explored all the different ways of removing the block from the hole and replacing it. You could just put it back in the same position. You could rotate it counterclockwise by 120 degrees and put it back in the hole. You could rotate it counterclockwise by 240 degrees and put it back. You could flip it over leaving the top corner fixed. You could flip it over leaving the top corner fixed and then rotate it counterclockwise by 120 degrees or by 240 degrees. You could combine two motions of the block by first doing one and then doing another, that is, you could compose one transformation with another to form a new transformation.



Now let's undertake the mathematical exploration of seeking the essential and isolating common features of these examples. All three examples involve combining two objects to get a third. In the case of addition of integers, we add two integers to get another integer;  $(2 + 3 \mapsto 5)$ . In the case of multiplication of reals, we multiply two reals to get another real number;  $(3 \cdot 1.204 \mapsto 3.612)$ . In the case of ways to move the block, we combine two transformations of the block to get a third;  $([\text{flip it over leaving the top corner fixed}] \circ [\text{flip it over fixing the top corner and rotate it by } 120^\circ \text{ counterclockwise}] \mapsto [\text{rotate it by } 240^\circ \text{ counterclockwise}])$ . Traditionally, if  $S$  and  $T$  are two transformations, then  $T \circ S$  means to perform the transformation  $S$  *then* the transformation  $T$ . If you read the symbol “ $\circ$ ” as “after”, you will do the transformations in the right order. Be careful to respect this convention.

What all of these examples have in common is that in each case we start with some collection (integers, reals, transformations of a triangle) and we have some operation (addition, multiplication, composition) that takes any two items from the collection and returns a third. Because our operations take *two* objects as input, we call them *binary* operations. Our rules for combining have some other features in common as well.

First common feature—an identity element: In each of our examples

there is an element that, when combined with any other element, has no effect on the other element. We call that “ineffective” element an *identity element*.

- (1) In addition of integers:  $0 + 3 = 3$ . In fact,  $0 + n = n = n + 0$  for every integer  $n$ .
- (2) In multiplication of reals:  $1 \cdot 2.35 = 2.35$ . In fact,  $1 \cdot r = r = r \cdot 1$  for every real number  $r$ .
- (3) In composing transformations of a triangular block: [just put the block back in the same position]  $\circ$  [rotate counterclockwise by  $120^\circ$ ] = [rotate counterclockwise by  $120^\circ$ ]. In fact, [just put the block back in the same position]  $\circ T = T = T \circ$  [just put the block back in the same position] for any transformation  $T$  of the block.

Second common feature—inverses: In each example, every element can be combined with another element to produce the identity element; that is, for each element there is another that undoes it. This “reversing” element is called an *inverse*.

Actually, in our example of the reals not every element has an inverse because nothing times 0 gives 1. So we will change the example of the reals under multiplication a little, namely, we will omit 0. Our second example will now be all the real numbers except 0. This process of modifying our examples in the face of difficulties has lead to lots of interesting mathematics.

- (1) In addition of integers:  $3 + (-3) = 0$ . In fact, for every integer  $n$ ,  $n + (-n) = 0 = (-n) + n$ .
- (2) In multiplication of reals except 0:  $2.35 \cdot \frac{1}{2.35} = 1$ . In fact, for every non-0 real number  $r$ ,  $r \cdot \frac{1}{r} = 1 = \frac{1}{r} \cdot r$ .
- (3) In composing transformations of a triangular block: [rotate counterclockwise by  $120^\circ$ ]  $\circ$  [rotate counterclockwise by  $240^\circ$ ] = [just put the block back in the same position]. In fact, every transformation of the block can be followed with another transformation that returns the block to its original position.

*Exercise 3.1.* Make a chart that lists each of the six transformations of the equilateral triangle and, for each transformation, find its inverse.

Third common feature—associativity: The rules for combining a pair of elements to get a third leaves us with an intriguing ambiguity about how three elements might be combined. When adding three integers, what do we do? Stop and compute  $2 + 4 + 6$ ; try to explain what you did. You probably reduced the question to a problem you knew how to deal with: add two of the integers and then add the result to the third. Similarly, when adding any number of integers, we start by adding two of them at a time and reduce until we get an answer.

So what does the expression  $a + b + c$  really mean? There are two different ways to break this expression down into a sequence of pair-wise addition problems:  $(a + b) + c$  or  $a + (b + c)$ . Parentheses mean what they always have, namely, the order of operations goes from inside the parentheses to outside. Both of these possible sequences are reasonable ways of reducing a question of adding three integers down to the case of adding pairs of integers sequentially. In our examples, the choice of sequencing doesn't matter. More precisely, in each example, both choices of ways to put parentheses on  $a + b + c$ ,  $x \cdot y \cdot z$ , or  $R \circ S \circ T$  produce the same result. This feature of the operation is called *associativity*.

- (1) In addition of integers: for any three integers  $a$ ,  $b$ , and  $c$ ,

$$(a + b) + c = a + (b + c).$$

- (2) In multiplication of reals except 0: for any three non-0 reals  $a$ ,  $b$ , and  $c$ ,

$$(a \cdot b) \cdot c = a \cdot (b \cdot c).$$

- (3) In composing transformations of a triangular block: for any three transformations  $R$ ,  $S$ , and  $T$ ,

$$(T \circ S) \circ R = T \circ (S \circ R).$$

This fact is not completely obvious, so you might fear that you'd have to verify it by laboriously checking every possible sequence of three transformations. Fortunately (for you and the grader), there is an easier way. These transformations are functions, and it is straight forward to show that the composition of functions is associative, as long as it is defined. This also explains the order convention of  $T \circ S$  as  $T$  after  $S$ .

*Exercise 3.2.* In composing transformations, check an example of associativity by confirming the following equality: ([flip fixing the top corner]  $\circ$



$[\text{rotate counterclockwise by } 120^\circ] \circ [\text{flip fixing the top corner then rotate counterclockwise by } 240^\circ] = [\text{flip fixing the top corner}] \circ ([\text{rotate counterclockwise by } 120^\circ] \circ [\text{flip fixing the top corner then rotate counterclockwise by } 240^\circ])$ .

*Exercise 3.3.* Let  $a, b, c, d$  be integers and consider the expression  $a+b+c+d$ . How many different ways are there to put parentheses on this expression such that only two integers are added at a time?

Let's note one feature that is not shared by all three of our examples. In the example of the integers under addition, for any integers  $a$  and  $b$ ,  $a+b = b+a$ . Likewise, in the example of the reals under multiplication, for any real numbers  $r$  and  $s$ ,  $r \cdot s = s \cdot r$ ; however, notice that the order makes a difference in composing transformations of the triangle.

*Exercise 3.4.* Find some examples of two transformations of an equilateral triangle where composing the transformations in one order gives a different result from doing them in the other order. Each of your examples should be a pair of transformations of the triangle,  $S$  and  $T$ , such that  $S \circ T \neq T \circ S$ .

When the order does not matter, that is, when we always get the same result no matter in which order we do the binary operation, then we call the operation **commutative**. We will talk more about this distinction later, but from Exercise 3.4 we know that it is possible that the same two elements combined in the opposite order might yield a different result.

Now let's take a step that creates mathematical ideas, namely, defining a concept that captures the common features that we have found. It turns out that we have isolated the essential ingredients of a mathematical structure that is called a *group*. We'll give the definition here and then make sure that we have pinned down all the features thoroughly.

*Definition.* A **group** is a set  $G$  with a binary operation  $*$ , written  $(G, *)$ , such that:

1. The operation  $*$  is closed and well-defined.
2. The operation  $*$  is associative.
3. There is an element  $e \in G$  such that  $g * e = g = e * g$  for all  $g \in G$ . The element  $e$  is called the **identity**. In particular  $G$  is non-empty.
4. For each element  $g \in G$  there is an element  $h \in G$  such that  $g * h = e = h * g$ . This element  $h$  is called the **inverse of  $g$**  and is often written as  $g^{-1}$ .

A binary operation is a procedure that takes two elements from a set and returns a third object. It is possible that this third object does not lie in our original set; if this happens, we say that the binary operation is not *closed*. Also, if the binary operation is given in terms of a rule, or if there is some ambiguity in the set, then sometimes the operation gives different values even though the input has not changed. If this happens, we say that the operation is not *well-defined*.

Our examples have given us an intuitive idea of what we want to convey, but we may want to take a further step of precision. In the Appendix *Sets and Functions* we clarify what we mean by a *set*, by a *function*, and by a *binary operation*. Here is an exercise to help you clarify the idea of a binary operation.

*Exercise 3.5.* Show that the following are or are not closed, well-defined binary operations on the given sets.

1. The interval  $[0, 1]$  with  $a * b = \min\{a, b\}$
2.  $\mathbb{R}$  with  $a * b = a/b$
3.  $\mathbb{Z}$  with  $a * b = a^2 + b^2$
4.  $\mathbb{Q}$  with  $a * b = \frac{\text{numerator of } a}{\text{denominator of } b}$
5.  $\mathbb{N}$  with  $a * b = a - b$

A group is a set with a closed, well-defined, associative binary operation  $(G, *)$  with an identity element,  $e$ , and an inverse element  $g^{-1}$  for each  $g \in G$ .

Let's begin our exploration of this new mathematical entity, a group, by first recording that our generative examples are groups. There is no need to prove these theorems now.

*Theorem 3.6.* The integers with addition,  $(\mathbb{Z}, +)$ , is a group.

*Theorem 3.7.* The non-zero real numbers with multiplication,  $(\mathbb{R} \setminus \{0\}, \cdot)$  is a group.

*Theorem 3.8.* The transformations of an equilateral triangle in the plane with composition is a group. We call this group  $D_3$ , the symmetries of the equilateral triangle.

When we write theorems about an arbitrary group  $(G, *)$ , we will often write  $G$  for  $(G, *)$  to simplify the notation; we know that  $G$  has a binary operation, but we don't explicitly name it. Similarly, we will sometimes write  $gh$  when we mean  $g * h$ . Our first theorem that is true for any group tells us that a group can have only one identity element.

*Theorem 3.9.* Let  $G$  be a group. There is a *unique* identity element in  $G$ . In other words, there is only one element in  $G$ ,  $e$ , such that  $g * e = e * g = g$  for all  $g$  in  $G$ .

Every group satisfies the following Cancellation Law. It seems simple and obvious, but it is an extremely useful property; it will reappear in the proof of every important theorem for the duration of this chapter.

*Theorem 3.10 (Cancellation Law).* Let  $G$  be a group, and let  $a, x, y \in G$ . Then  $a * x = a * y$  if and only if  $x = y$ .

Be careful not to use the theorem when proving it. Instead, use associativity and the well-definedness of  $*$ . As always, the phrase “if and only if” means that there are actually two theorems here to prove:  $a * x = a * y$  implies  $x = y$ , and  $x = y$  implies  $a * x = a * y$ .

*Exercise 3.11.* Show that the Cancellation Law fails for  $(\mathbb{R}, \cdot)$ , thus confirming that  $(\mathbb{R}, \cdot)$  is not a group.

*Corollary 3.12.* Let  $G$  be a group. Then each element  $g$  in  $G$  has a *unique* inverse in  $G$ . In other words, for a fixed  $g$ , there is only one element,  $h$ , such that  $g * h = h * g = e$ .

Recall that in general  $g * h$  may not equal  $h * g$ ; however, if one product is the identity, then both orders of the product yield the identity.

*Theorem 3.13.* Let  $G$  be a group with elements  $g$  and  $h$ . If  $g * h = e$ , then  $h * g = e$ .

In words, this theorem says that, if  $h$  is a right inverse of  $g$ , then it is also a left inverse of  $g$ . So we only need to check that  $g$  and  $h$  are one-sided inverses to know that they are inverses.

*Theorem 3.14.* Let  $G$  be a group and  $g \in G$ . Then  $(g^{-1})^{-1} = g$ .

Theorems like the preceding four show us that if we have a structure that satisfies the definition of a group, then it will automatically have the features stated in the theorems. One of the strategies and strengths of abstract mathematics is that we define a structure (like a group) and then deduce that any mathematical object of that type (any group, for example) will have features (like a unique identity or the cancellation property) that are common to every example of such a structure (every group).

In order to develop our intuition about groups, let's first consider a few more examples that we can create by taking our existing examples and seeking variations. Taking examples and concepts that we have and making variations of them is one of the most common and most powerful methods for creating new mathematical ideas.

Our first example of a group was the integers with the binary operation of addition,  $(\mathbb{Z}, +)$ . In life we also perform addition of numbers when we tell time, but in that case we have a cyclical kind of addition. If it's 9 o'clock now, then in 47 hours it will be 8 o'clock. Somehow in our world of time, " $9 + 47 = 8$ ". Can we construct a group that captures this cyclical kind of arithmetic? Well, we know what we need to construct a group: we need a set of elements and a binary operation. So to construct a group that captures the idea of times of the day, we might consider the hours  $\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$  as our set and *clock addition* as our operation. Notice that this definition of clock addition only allows us to combine two numbers from 1 to 12, so we could not add 47 to a time, for example. We'll deal with that issue later. For now, we have created a group.

*Exercise 3.15.* If  $G = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$  and  $\oplus$  is the binary operation clock addition, show that  $(G, \oplus)$  is a group. What is the identity element? What is the inverse of 3?

Once we have defined this clock group, we cannot (and should not) resist the urge to extend the idea. An obvious and important way to generalize the idea is to consider clock arithmetic with different numbers of hours in the day. That generalization gives us infinitely many different groups that use cyclical addition. Let's now take the step of pinning down all these ideas with formal definitions.

Let  $C_n = \{0, 1, \dots, n-1\}$ . We will define a binary operation on  $C_n$  that captures the idea of cyclical arithmetic. Any two elements in  $C_n$  are integers, so we can add them. If their sum is less than  $n$ , then it is in  $C_n$ , so their sum makes sense as an element of  $C_n$ . If their sum is  $m$ , bigger than  $n$ , then replace it by  $m - n$ , which is now definitely back in the set. Call this operation *n-cyclic addition*, and write it as  $\oplus_n$ .

*Theorem 3.16.* For every natural number  $n$ , the set  $C_n$  with *n-cyclic addition*,  $(C_n, \oplus_n)$ , is a group. We call it the **cyclic group of order  $n$** .

Our cyclic groups are nice, but somehow we need to deal with the fact that in reality we can add 47 hours to a time. How can we extend our description of our clock world so as to include 47 and other integers in it? A solution is presented to us by considering European and military time, where time is measured with a 24 hour clock rather than a 12 hour clock. When their clocks read 15 o'clock, ours read 3 o'clock. This idea of reducing by 12 can easily be extended even to 47. What time is 47 o'clock? Answer: it's 11, because  $47 - 12$  is 35,  $35 - 12$  is 23, and  $23 - 12$  is 11. More simply put, since  $47 = 3 \cdot 12 + 11$  or, equivalently,  $47 - 11 = 3 \cdot 12$ , we consider 47 and 11 to be referring to the same time, that is to say, 47 and 11 should

be different names for the same element of our group that captures the idea of time. In general, we could say that two integers  $a$  and  $b$  are equivalent if  $a = b + 12k$  for some integer  $k$  or, equivalently,  $a - b$  is a multiple of 12. With this idea in mind, we can think of a new group with twelve elements  $\{[1], [2], [3], \dots, [12]\}$ ; however, each element really stands for all the integers that are equivalent to it using our concept of time equality.

*Definition.* Let  $n$  be a natural number. Two integers  $a$  and  $b$  are said to be **congruent modulo  $n$**  if there exists another integer  $k$  such that  $a = b + kn$  or, equivalently,  $a - b = kn$ . In other words, two integers are congruent modulo  $n$  if and only if their difference is divisible by  $n$ . We write “ $a$  is congruent to  $b$  modulo  $n$ ” as  $a \equiv b \pmod{n}$ .

*Definition.* Let  $\mathbb{Z}_n = \{[a]_n \mid a \in \mathbb{Z}, [a]_n = [b]_n \text{ if and only if } a \equiv b \pmod{n}\}$ . Then define the binary operation  $\oplus$  on  $\mathbb{Z}_n$  by  $[a]_n \oplus [b]_n = [a + b]_n$ , which we will call **modular addition**.

Now we can now try answering the question of “what time is it 47 hours after 8 o’clock?” Using modular arithmetic, we can replace the question of  $9 \oplus 47$  with  $[9]_{12} \oplus [47]_{12} = [56]_{12} = [8 + 4(12)]_{12} = [8]_{12}$ .

*Exercise 3.17.* Show that  $\oplus$  is well-defined on  $\mathbb{Z}_n$ . That means, show that if you replace integers  $a$  and  $b$  by congruent integers  $a'$  and  $b'$  respectively, then  $[a]_n \oplus [b]_n = [a']_n \oplus [b']_n$ .

Note that both  $(C_n, \oplus_n)$  and  $(\mathbb{Z}_n, \oplus)$  are groups with  $n$  elements and a cyclical addition. Intuitively, they are clearly the “same”, but we do not yet have a definition for when two groups are the same. We will return to this issue later, but for now, notice that when working with  $C_n$ , the group elements look like  $\{0, 1, \dots, n-1\}$ , and the operation is called  $n$ -cyclic addition, written  $\oplus_n$ . When working with  $\mathbb{Z}_n$ , the distinct elements look like  $\{[0]_n, [1]_n, \dots, [n-1]_n\}$ , and the operation is called modular addition, written  $\oplus$ .

These families of groups were suggested as variations on the group of integers with addition. Let’s now turn to the task of generalizing and extending another one of our generative examples, the symmetries of the triangle.

## 3.2 Symmetry Groups of Regular Polygons

The symmetries of an equilateral triangle under composition form one of our generative examples of a group. We called this group  $D_3$  because it consisted of transformations of a shape in the plane with 3 sides, having uniform side lengths and angles. Similarly, we could think about the transformations of

the shape with 4 equal sides and uniform angles, more commonly known as the square. Of course, we could also consider similar shapes with any number of sides. So we can create related groups by considering the symmetries of any *regular polygon*.

By a symmetry, we mean a transformation that takes the polygon to itself as a rigid object, which is to say that we consider only rotations and reflections. The elements of these groups of symmetries are functions, so we can use composition as the binary operation.

*Exercise 3.18.* Consider a square in the plane. How many distinct symmetries does it have? Give each symmetry an intuitive label. For every pair of symmetries,  $S_1$  and  $S_2$ , compose them in both orders,  $S_2 \circ S_1$  and  $S_1 \circ S_2$ . That is, for each composition, figure out which label from your list it equals. Try to find a convenient way to record all of this information.

The following theorem notes that we have a new group.

*Theorem 3.19.* The symmetries of the square in the plane with composition form a group.

The symmetries of the regular 4-gon (i.e., the square) with the binary operation of composition is denoted  $D_4$ . In general, the symmetries of the regular  $n$ -gon form a group, which we call  $D_n$ .

*Exercise 3.20.* For each natural number  $n$ , how many elements does  $D_n$  have?

For groups with a finite number of elements, it is possible to write out a table of all group elements and the result of the binary operation on any pair. Such a table is called a Cayley table. For example, the following chart is the Cayley table for  $(C_5, \oplus_5)$ .

$\oplus_5$	0	1	2	3	4
0	0	1	2	3	4
1	1	2	3	4	0
2	2	3	4	0	1
3	3	4	0	1	2
4	4	0	1	2	3

*Exercise 3.21.* Write out a Cayley table for  $D_4$ .

### 3.3 Subgroups, Generators, and Cyclic Groups

A **subgroup** of a group  $(G, *)$  is a non-empty subset  $H$  of  $G$  along with the restricted binary operation such that  $(H, *|_H)$  is a group. Checking that a

subset is a subgroup is just like checking the axioms of a group, though more attention is often paid to the subset being closed under  $*$ ; that is, we need to check that for any two elements in the subgroup, the binary operation performed on them results in another element of the subgroup.

Let's take a look at our examples and find some of their subgroups.

- Exercise 3.22.*
1. Show that the even integers, written  $2\mathbb{Z}$ , form a subgroup of  $(\mathbb{Z}, +)$ . Technically,  $2\mathbb{Z}$  is just a set, but we will often drop the binary operation from the notation for a subgroup when it is obvious.
  2. Show that the set of non-zero rational numbers,  $\mathbb{Q} \setminus \{0\}$ , is a subgroup of  $(\mathbb{R} \setminus \{0\}, \cdot)$ .
  3. Show that the set of three transformations  $H = \{[\text{just put the triangle back where you found it}], [\text{rotate counterclockwise by } 120^\circ], [\text{rotate counterclockwise by } 240^\circ]\}$  is a subgroup of  $D_3$ , the symmetries of the triangle.
  4. Show that  $K = \{0, 15, 30, 45\}$  is a subgroup of  $(C_{60}, \oplus_{60})$ .

If a group is defined by a set of objects satisfying a certain condition, then a subgroup is usually a subset satisfying a stronger condition. In the exercise above, the *even integers* form a subgroup of *all integers*, the non-zero *rational* numbers are a subgroup of the non-zero *reals*, the *rotations* are a subgroup of *all symmetries*, and the *quarter hours* are a subgroup of *all minutes* in a group capturing 60-minute time arithmetic.

But not every condition defines a subgroup. For example, the odd integers are not a subgroup of the integers under addition. (Why not?) To get a sense of what subsets of a group form a subgroup, it is a good exercise to describe all the subgroups of a few groups.

*Exercise 3.23.* For each of the following groups, find all subgroups. Argue that your list is complete.

1.  $(D_4, \circ)$
2.  $(\mathbb{Z}, +)$
3.  $(C_n, \oplus_n)$

You may have noticed that the identity element is in each of your subgroups.

*Theorem 3.24.* Let  $G$  be a group with identity element  $e$ . Then for every subgroup  $H$  of  $G$ ,  $e \in H$ .

The smallest and simplest subgroup of any group is just the identity element.

*Theorem 3.25.* Let  $G$  be a group with identity element  $e$ . Then  $\{e\}$  is a subgroup of  $G$ .

Every group is a subgroup of itself; this subgroup is necessarily the biggest subgroup.

*Theorem 3.26.* Let  $G$  be a group. Then  $G$  is a subgroup of  $G$ .

Any group  $G$  has the subgroups  $\{e\}$  and  $G$ , so these subgroups are basically trivial. If  $H$  is a subgroup of  $G$  such that  $\{e\} \subsetneq H \subsetneq G$ , then we say that  $H$  is a **non-trivial** subgroup of  $G$ .

For repeated applications of the binary operation to one element, we will sometimes use exponents:  $g^4 = g * g * g * g$ . We have some intuition about how exponents work, and that intuition is the reason for this shorthand notation, but be careful not to use any properties of exponents that you have not checked for groups. If  $n$  is a positive integer, then  $g^n$  is the product of  $n$  copies of  $g$ . If  $n$  is a negative integer, then  $g^n$  is the product of  $-n$  copies of  $g^{-1}$ . We have already defined  $g^{-1}$  as the inverse of  $g$ , and we now define  $g^1 = g$ . Also, we let  $g^0 = e$ .

*Exercise 3.27.* Let  $G$  be a group,  $g \in G$ , and  $n, m \in \mathbb{Z}$ . Then

1.  $g^n g^m = g^{n+m}$ , and
2.  $(g^n)^{-1} = g^{-n}$ .

*Definition.* Let  $G$  be a group and  $g$  be an element of  $G$ . Then  $\langle g \rangle$  is the subset of elements of  $G$  formed by repeated applications of the binary operation using only  $g$  and  $g^{-1}$ , that is,  $\langle g \rangle = \{g^{\pm 1} * g^{\pm 1} * \dots * g^{\pm 1}\}$ . Using the notation of the previous exercise,  $\langle g \rangle = \{g^m \mid \text{for all } m \in \mathbb{Z}\}$ .

*Theorem 3.28.* Let  $G$  be a group and  $g$  be an element of  $G$ . Then  $\langle g \rangle$  is a subgroup of  $G$ .

We call  $\langle g \rangle$  the **subgroup of  $G$  generated by  $g$** . Notice that the even integers can all be thought of as repeated additions of 2 and  $-2$ . In other words,  $2\mathbb{Z}$  is the subgroup of  $(\mathbb{Z}, +)$  generated by 2 (or  $-2$ ), that is,  $2\mathbb{Z} = \langle 2 \rangle = \langle -2 \rangle$ .

Similarly we can talk about subgroups generated by more than one element. If  $S$  is any subset of a group  $G$ , then we define  $\langle S \rangle$  to be all elements of  $G$  that are obtained from finite combinations of elements of  $S$  and their inverses,  $\{s_1^{\pm 1} * s_2^{\pm 1} * \dots * s_n^{\pm 1} \mid s_i \in S\}$ .



*Theorem 3.29.* Let  $G$  be a group and  $S$  be a subset of  $G$ . Then  $\langle S \rangle$  is a subgroup of  $G$ .

As was the case with a single element,  $\langle S \rangle$  is called the subgroup generated by  $S$ . It is the smallest subgroup that contains all the elements of  $S$ , and  $\langle g \rangle$  is the smallest subgroup that contains  $g$ .

*Theorem 3.30.* Let  $G$  be a group,  $S$  be a subset of  $G$ , and  $H$  be a subgroup of  $G$  such that  $S \subset H$ . Then  $\langle S \rangle$  is a subgroup of  $H$ .

The preceding theorems show us a method for constructing a subgroup of any group. We can just start with any collection of elements from the group and then look at all the elements we get by performing the binary operation repeatedly on those elements and their inverses.

*Exercise 3.31.* Which subgroup of  $\mathbb{Z}$  is  $\langle 5, -8 \rangle$ ?

*Definition.* A group  $G$  is called **cyclic** if there is an element  $g$  in  $G$  such that  $\langle g \rangle = G$ . In other words, a group is cyclic if it is generated by one element.

Some of the groups that we have considered are cyclic.

*Theorem 3.32.* The integers under addition is a cyclic group.

*Theorem 3.33.* For every natural number  $n$ , the groups  $C_n$  and  $\mathbb{Z}_n$  are cyclic groups.

In general, subgroups of groups can be complicated; however, subgroups of cyclic groups are all generated by one element.

*Theorem 3.34.* Any subgroup of a cyclic group is cyclic.

*Exercise 3.35.* For each natural number  $n$ , what is the smallest number of elements that generate  $D_n$ ?

*Definition.* A group  $G$  is called **finite** if the underlying set is finite. Similarly,  $G$  is called **infinite** if its underlying set is infinite. A group is **finitely generated** if  $G = \langle S \rangle$  for some finite subset  $S$  of its elements.

*Theorem 3.36.* Every finite group  $G$  is finitely generated.

*Theorem 3.37.* The groups  $(\mathbb{R} \setminus \{0\}, \cdot)$  and  $(\mathbb{Q} \setminus \{0\}, \cdot)$  are not finitely generated.

*Definition.* The number of elements in (the underlying set of)  $G$  is called the **order** of  $G$ , written  $|G|$ . The order of an element  $g$ , written  $o(g)$ , is the order of the subgroup that it generates,  $o(g) = |\langle g \rangle|$ .

*Exercise 3.38.* Compute the order of each element  $T \in D_4$ . Carefully use the definition of  $o(T)$ .

The order of an element of a group is defined in terms of the number of elements in  $\langle g \rangle$ ; however, that number is also the smallest power of the element that equals the identity element of the group.

*Theorem 3.39.* Let  $g$  be an element of a finite group  $G$  whose identity element is  $e$ . Then  $o(g) = |\langle g \rangle|$  is the smallest natural number  $r$  such that  $g^r = e$ .

You may have noticed a fundamental difference between the structures of  $(\mathbb{Z}, +)$  and  $(D_n, \circ)$ . For example, when adding integers, the order doesn't matter;  $a + b = b + a$  for any pair of integers. But this order does matter quite a bit when composing functions/symmetries. We give a special name to groups whose operation is commutative.

*Definition.* A group  $(G, *)$  is **abelian** if and only if, for every pair of elements  $g, h \in G$ ,  $g * h = h * g$ . So, a group is abelian if and only if its binary operation is commutative.

All cyclic groups are abelian.

*Theorem 3.40.* If  $G$  is a cyclic group, then  $G$  is abelian.

The commutativity of integer addition seems like a fundamental assumption about mathematics, since it is something we each learned at such a young age. But we can actually prove this fact just using the properties of  $\mathbb{Z}$  as a group, since all cyclic groups are abelian.

*Corollary 3.41.* The integers under addition form an abelian group.

Not all groups are cyclic, and we have already met a few non-cyclic groups.

*Corollary 3.42.* The group  $D_4$  is not cyclic.

Although cyclic groups are abelian, there are abelian groups that are not cyclic.

*Exercise 3.43.* Give an example of a finite group that is abelian but not cyclic. The smallest such group has four elements and is most easily described by writing its Cayley table.

In abelian groups every element commutes with every other element. In non-abelian groups, there can still be some elements that commute with all of the elements. We know that the identity element always commutes with every element, for example. We will name the set of elements that commute with every element of a group.

*Definition.* The **center** of a group  $G$  is the collection of elements in  $G$  that commute with all of the elements. The center is denoted  $Z(G)$  and can be described as

$$Z(G) = \{g \in G \mid g * h = h * g \text{ for all } h \in G\}.$$

The center of a group is not just a collection of elements of the group, it is a subgroup of the group.

*Theorem 3.44.* Let  $G$  be a group. Then  $Z(G)$  is a subgroup of  $G$ .

*Exercise 3.45.* Give examples of groups  $G$  in which

1.  $Z(G) = \{e\}$ ;
2.  $Z(G) = G$ ; and
3.  $\{e\} \subsetneq Z(G) \subsetneq G$ .

### 3.4 Products of Groups

We've described a few interesting groups, and we know that one way to find other groups is to find subgroups of the ones we have. Another method for building new sets from existing ones is to take their Cartesian product. If  $A$  and  $B$  are sets, then  $A \times B = \{(a, b) | a \in A \text{ and } b \in B\}$ , the set of ordered pairs of elements from  $A$  and  $B$ , which is called their *Cartesian Product*. We can make the Cartesian product of two groups into a group.

*Theorem 3.46.* Let  $(G, *_G)$  and  $(H, *_H)$  be groups and define  $*$  :  $(G \times H) \times (G \times H) \rightarrow G \times H$  by  $(g_1, h_1) * (g_2, h_2) = (g_1 *_G g_2, h_1 *_H h_2)$ . Then  $(G \times H, *)$  is a group, called the **(direct) product** of  $G$  and  $H$ .

The next exercise asks you to explore when the direct product of two cyclic groups is or is not a cyclic group.

*Exercise 3.47.* For natural numbers  $n$  and  $m$ , when is the group  $\mathbb{Z}_n \times \mathbb{Z}_m$  cyclic?

The direct product of cyclic groups may not always be cyclic; however, the direct product of abelian groups is always abelian.

*Theorem 3.48.* Let  $G$  and  $H$  be groups. Then  $G \times H$  is abelian if and only if both  $G$  and  $H$  are abelian.

The direct product of two groups has some natural subgroups.

*Theorem 3.49.* Let  $G_1$  be a subgroup of a group  $G$  and  $H_1$  be a subgroup of a group  $H$ . Then  $G_1 \times H_1$  is a subgroup of  $G \times H$ .

If we can realize a complicated group as the direct product of smaller groups, then we can feel that we know a lot about its structure. One of the most famous such structure theorems tell us that every finite abelian group is the direct product of cyclic groups.

### 3.5 Symmetric Groups

Many interesting groups have elements each of which is a function. We've already considered some such groups, so let's start there. The function groups we have encountered thus far are the symmetry groups of the regular  $n$ -gons,  $D_n$ . The elements are rigid transformations, which are functions from the  $n$ -gon to itself, and the operation is composition.

But what properties of these functions were necessary and what were superfluous for constructing a group of functions? Let's think about what properties of functions relate to the four properties of a group. First, the fact that the functions in  $D_n$  were from a set to itself made composition make sense in all cases, so the binary operation will be closed and well-defined. Associativity is automatic for composition of functions, so we won't need to worry about that. You may recall that a function needs to be injective to have an inverse under composition, and for that inverse to have the correct domain, the original function needs to be surjective. Putting these observations together, we should try to make a group whose elements are bijective functions from a set to itself. (If these observations are not familiar, see the Appendix *Sets and Functions*.)

*Theorem 3.50.* Let  $X$  be a set, let  $Sym(X)$  be the set of bijections from  $X$  to  $X$ , and let  $\circ$  represent composition. Then  $(Sym(X), \circ)$  is a group.

In the case where  $X$  is a finite set with  $n$  elements, we will usually write  $S_n$  for  $Sym(X)$ .

*Theorem 3.51.* If  $X$  is a set with  $n$  elements, then  $|Sym(X)| = |S_n| = n!$ .

The group  $S_n$  is called the **symmetric group on  $n$  points** or objects. Label the  $n$  elements of the set  $X$  as  $\{1, 2, \dots, n\}$ . The elements of  $S_n$  are bijective functions, which is to say that each element of  $S_n$  is just a permutation of these  $n$  points. For example, in  $S_6$  one element is the function from  $\{1, 2, 3, 4, 5, 6\}$  to  $\{1, 2, 3, 4, 5, 6\}$  defined by  $1 \mapsto 3$ ,  $2 \mapsto 5$ ,  $3 \mapsto 2$ ,  $4 \mapsto 6$ ,  $5 \mapsto 1$ , and  $6 \mapsto 4$ ; call it  $g$ . We can denote this permutation in a clever way:

$$g = (1325)(46).$$

Here is how to interpret the notation: each element is sent to the element to its right until we reach a close parenthesis symbol, ')', in which case the element right before the close parenthesis gets sent to the first element within that pair of parentheses. To generate this notation, we started by drawing an open parenthesis and writing 1; next to it we write where 1 is sent by  $g$ , namely 3. Next to that we write where 3 is sent, namely, 2. We continue

doing this until a number, in this case 5, is sent to one that we've already written (which will be the first one we wrote next to the open parenthesis) at which point we close the parentheses. If there are any unused numbers, we repeat the process with a new open parenthesis symbol '(' and one of the unused elements; in this case, we start with 4 and continue to 6 and then close the parentheses when we find that 6 returns to 4. If there were an element that was sent to itself by the permutation, for example  $7 \mapsto 7$ , it would be represented as (7) in the parenthesis notation.

*Exercise 3.52.* Let  $h$  be the element of  $S_8$  that sends  $1 \mapsto 4$ ,  $2 \mapsto 7$ ,  $3 \mapsto 8$ ,  $4 \mapsto 6$ ,  $5 \mapsto 5$ ,  $6 \mapsto 1$ ,  $7 \mapsto 2$ , and  $8 \mapsto 3$ . Write down the parenthesis notation for  $h$ .

*Exercise 3.53.* Suppose you are given the parenthesis notation for two symmetries of  $S_n$ . If you are given the two elements of  $S_n$ ,  $g$  and  $h$ , what is the parenthesis notation for the composition  $h \circ g$ ? Recall that we're thinking about  $g$  and  $h$  as functions, so  $h \circ g$  means do the permutation  $g$  first and then  $h$ . The point of this exercise is for you to figure out and describe how the composition of permutations can be executed using the notation. In particular, if  $g = (132)(456)$  and  $h = (1463)(25)$ , compute  $h \circ g$ .

Reading off the numbers in one set of parentheses gives a sort of circuit that an element would follow under repeated applications of that particular permutation.

*Exercise 3.54.* Suppose  $g \in S_n$  and that you know the parenthesis notation for  $g$ . How can you compute  $o(g)$  without repeatedly composing  $g$  with itself?

*Exercise 3.55.* Write out the parenthesis notation for all elements of  $S_4$ .

*Exercise 3.56.* Find all subgroups of  $S_4$ .

As we mentioned above, the groups  $D_n$  are also bijections from a set to itself. Within the group of bijections, these are the rigid transformations. This means that these groups are subgroups of symmetric groups.

*Exercise 3.57.* Rewrite  $D_4$  as a subgroup of  $S_4$  by thinking of each symmetry in  $D_4$  as a permutation of the 4 corners of the square.

It turns out that *any* group can be thought of as a subgroup of a symmetric group. To prove this fact, our challenge is to associate an arbitrary element of an arbitrary group with a bijection on some set.

*Exercise 3.58.* Let  $G$  be a group. For each element  $g \in G$ , define the function  $\lambda_g : G \rightarrow G$  by  $\lambda_g(h) = gh$ . Check that, for any  $g \in G$ ,  $\lambda_g$  is a bijection and that  $\lambda_g \neq \lambda_{g'}$  when  $g \neq g'$ .

*Exercise 3.59.* Let  $G$  be a group. Let  $\Lambda : G \rightarrow \text{Sym}(G)$  be the function  $\Lambda(g) = \lambda_g$ . Show that  $\Lambda$  is injective.

One of the strategies of mathematical exploration is to find the most general or most comprehensive examples of a mathematical object. The previous exercise suggests that understanding symmetric groups and their subgroups amounts to understanding *all* groups. Unfortunately, another way to look at these insights is to say that the symmetric groups are as complicated as any groups that exist, so they will be difficult to fully fathom. In any case, thinking about elements of groups as permutations is often a valuable strategy. We'll talk more about this later, in the Section *Groups in Action*.

### 3.6 Maps between Groups

After we have defined a mathematical object like a group, we should be able to define what it means for two such objects to be considered “the same”. Our concept of “sameness” should depend on what we view as the fundamental, defining features of the object in question. In the case of a group, the definition tells the story: a group is a non-empty set together with a binary operation. So if we look at two groups, we want the concept of “sameness” to refer to the sets involved and their respective binary operations. Pinning this idea down is a basic strategy for exploring a mathematical idea. Once we have defined a mathematical object (in this case a group), we can ask what kind of functions between these objects (groups) respect the defining structure of the object (the sets and binary operations). Let's see what this abstract philosophy means in the case of groups.

For two groups  $G$  and  $H$  to be the same, their underlying sets should be in bijective correspondence. Thinking in terms of finite groups, there should be a relabeling of the elements that makes the Cayley tables look identical. Computing the binary operation before or after the relabeling should not matter. When two groups are the same in this sense that one group is just a relabeling of the elements of the other, then we call the groups *isomorphic*.

*Exercise 3.60.* Using this informal definition of isomorphic, show that  $D_3$  and  $S_3$  are isomorphic.

The way we formalize the idea that two sets,  $X$  and  $Y$ , are relabelings of each other is by finding a bijection  $f : X \rightarrow Y$ . So for two groups  $G$  and  $H$  to be the same, there must be a bijection  $\phi : G \rightarrow H$ . But, in addition, the relabeling of the elements should respect the binary operations. Suppose that the elements  $a, b, c$  in the group  $G$  correspond respectively to

the elements  $A, B, C$  in the group  $H$ . Then  $a *_G b = c$  (in the group  $G$ ) should mean that  $A *_H B = C$  (in the group  $H$ ). If  $\phi : G \rightarrow H$  is the function that defines the relabeling, then we're saying  $\phi(a) = A$ ,  $\phi(b) = B$ , and  $\phi(c) = C$ . Written in this notation, the desired equality becomes

$$\phi(a *_G b) = \phi(c) = C = A *_H B = \phi(a) *_H \phi(b).$$

Now we're ready to formalize what it means for a function,  $\phi : (G, *_G) \rightarrow (H, *_H)$ , to “respect the binary operations”.

*Definition.* Let  $(G, *_G)$  and  $(H, *_H)$  be groups and let  $\phi : G \rightarrow H$  be a function on their underlying sets. Then we call  $\phi$  a **homomorphism** of the groups if for every pair of elements  $g_1, g_2 \in G$ ,  $\phi(g_1 *_G g_2) = \phi(g_1) *_H \phi(g_2)$ .

Notice that the central point in the definition of a homomorphism is that the binary operation in the domain is related to the binary operation in the range; first doing the binary operation between two elements in  $G$  and then performing the homomorphism gives the same element in  $H$  as first performing the homomorphism on the two elements individually and then combining the results with the binary operation in  $H$ . “Combine then map” should be the same as “map then combine”.

Let's begin by examining a few homomorphisms between pairs of our favorite groups.

*Exercise 3.61.* Confirm that each of the following functions is a homomorphism.

1.  $\phi : \mathbb{Z}_{12} \rightarrow \mathbb{Z}_{24}$  defined by  $\phi([a]_{12}) = [2a]_{24}$ . (However, notice that  $\phi([a]_{12}) = [a]_{24}$  is not a homomorphism.)
2.  $\phi : \mathbb{Z} \rightarrow \mathbb{Z}_9$  defined by  $\phi(a) = [3a]_9$
3.  $\phi : \mathbb{Z}_6 \rightarrow \mathbb{Z}_3$  defined by  $\phi([a]_6) = [a]_3$
4.  $\phi : \mathbb{Z}_n \rightarrow \mathbb{Z}_n$  defined by  $\phi([a]_n) = [-a]_n$

*Definition.* Let  $A$  be a subset of a set  $B$ . The **inclusion map**  $i_A : A \rightarrow B$  is defined as follows: for each element  $a \in A$ ,  $i_A(a) = a$ .

*Theorem 3.62.* Let  $H$  be a subgroup of a group  $G$ . Then the inclusion of  $H$  into  $G$ ,  $i_H : H \rightarrow G$ , is a homomorphism.

Let's see some basic consequences of the definition of homomorphisms. The next theorem tells us that any homomorphism takes the identity of the domain group to the identity of the codomain group. If there are several groups floating around, we may write  $e_G$  for the identity of  $G$ .

*Theorem 3.63.* If  $\phi : G \rightarrow H$  is a homomorphism, then  $\phi(e_G) = e_H$ .

Similarly, this next theorem tells us that homomorphisms send inverses to inverses.

*Theorem 3.64.* If  $\phi : G \rightarrow H$  is a homomorphism and  $g \in G$ , then  $\phi(g^{-1}) = [\phi(g)]^{-1}$ .

The next theorem tells us that the group structure is preserved by homomorphisms in the sense that the image of a group is a subgroup of the codomain.

*Definition.* For any function,  $f : A \rightarrow B$ , the **image** of  $f$  is

$$Im(f) = \{b \in B \mid \text{there exists an } a \in A \text{ such that } f(a) = b\},$$

which is also called the range of  $f$ . In words, the image of a function is the set of elements in the codomain that are “hit” by elements from the domain.

*Theorem 3.65.* If  $\phi$  is a homomorphism from the group  $G$  to the group  $H$ , then  $Im(\phi)$ , the image of  $\phi$ , is a subgroup of  $H$ .

*Theorem 3.66.* Let  $G$  and  $H$  be groups, let  $K$  be a subgroup of the group  $G$ , and let  $\phi : G \rightarrow H$  be a homomorphism, then  $\phi(K)$  is a subgroup of  $H$ .

*Exercise 3.67.* Let  $G$  be a group with an element  $g$ . Then define the function  $\phi : \mathbb{Z} \rightarrow G$  by setting  $\phi(n) = g^n$  for each  $n \in \mathbb{Z}$ . Show that  $\phi$  is a homomorphism. What is the image of  $\phi$ ?

*Exercise 3.68.* Let  $G = \langle g \rangle$  be a cyclic group and  $\phi : G \rightarrow H$  a homomorphism. Show that  $\phi(g)$  determines  $\phi(g')$  for all  $g' \in G$ .

The subgroup-preserving property of images of homomorphisms also works with the preimages of subgroups of the codomain under a homomorphism.

*Theorem 3.69.* Let  $G$  and  $H$  be groups, let  $L$  be a subgroup of the group  $H$ , and let  $\phi : G \rightarrow H$  be a homomorphism, then  $\phi^{-1}(L) = \{g \in G \mid \phi(g) \in L\}$  is a subgroup of  $G$ .

One such preimage is so important that it has a name of its own.

*Definition.* Let  $\phi : G \rightarrow H$  be a homomorphism from a group  $G$  to a group  $H$ . Then  $Ker(\phi) = \{g \in G \mid \phi(g) = e_H\}$  is called the **kernel** of  $\phi$ .

*Corollary 3.70.* For any homomorphism  $\phi : G \rightarrow H$ ,  $Ker(\phi)$  is a subgroup of  $G$ .



When we described functions in the Appendix *Sets and Functions*, we described several types of functions: injective (1-1) functions, surjective (onto) functions, and bijective (1-1 and onto) functions. In Exercise 3.61 above, we saw examples of homomorphisms which, as functions, fell into each of these categories. We give special names to each of these types of homomorphisms.

*Definitions.* 1. An injective homomorphism is called a **monomorphism**.

2. A surjective homomorphism is called an **epimorphism**.

3. A bijective homomorphism is called an **isomorphism**.

4. A group  $G$  is **isomorphic** to a group  $H$  if there is an isomorphism,  $\phi : G \rightarrow H$ , written  $G \cong H$ .

To get accustomed to these terms, let's begin by classifying each homomorphism from Exercise 3.61 above as a monomorphism, epimorphism, or isomorphism.

*Exercise 3.71.* Classify each of the following homomorphisms as a monomorphism, an epimorphism, an isomorphism, or none of these special types of homomorphisms.

1.  $\phi : \mathbb{Z}_{12} \rightarrow \mathbb{Z}_{24}$  defined by  $\phi([a]_{12}) = [2a]_{24}$ .
2.  $\phi : \mathbb{Z} \rightarrow \mathbb{Z}_9$  defined by  $\phi(a) = [3a]_9$
3.  $\phi : \mathbb{Z}_6 \rightarrow \mathbb{Z}_3$  defined by  $\phi([a]_6) = [a]_3$
4.  $\phi : \mathbb{Z}_n \rightarrow \mathbb{Z}_n$  defined by  $\phi([a]_n) = [-a]_n$

The integers can map onto any one of the modular arithmetic groups,  $\mathbb{Z}_n$ , by a homomorphism.

*Exercise 3.72.* For any natural number  $n$ , there is an epimorphism  $\phi_n : \mathbb{Z} \rightarrow \mathbb{Z}_n$ .

Projections from a direct product of groups to one of the factors are examples of epimorphisms. Let's define our terms.

*Definition.* Let  $X$  and  $Y$  be sets with Cartesian product  $X \times Y$ . The functions  $\pi_X((x, y)) = x$  and  $\pi_Y((x, y)) = y$  are called **projections** to the first and second coordinates respectively.

*Theorem 3.73.* Let  $G$  and  $H$  be groups. Then the projection maps  $\pi_G : G \times H \rightarrow G$  and  $\pi_H : G \times H \rightarrow H$  are epimorphisms.

The next theorem relates homomorphisms into groups with a homomorphism into their direct product.

*Theorem 3.74.* Let  $G$ ,  $H$ , and  $K$  be groups with homomorphisms  $f_1 : K \rightarrow G$  and  $f_2 : K \rightarrow H$ . Then there is homomorphism  $f : K \rightarrow G \times H$  such that  $\pi_G \circ f = f_1$  and  $\pi_H \circ f = f_2$ . Furthermore,  $f$  is the only function satisfying these properties. Moreover, if either  $f_1$  or  $f_2$  is a monomorphism, then  $f$  is also a monomorphism.

The concept of isomorphism is extremely important because two groups being isomorphic captures the idea that the two groups are the “same” by formalizing the notion that the two groups are just relabelings of each other. We have already been introduced to two groups that should be the same: the cyclic arithmetic and modular arithmetic groups of the same order. Let’s confirm this feeling by showing that they are isomorphic.

*Theorem 3.75.* The two groups  $(C_n, \oplus_n)$  and  $(\mathbb{Z}_n, \oplus)$  are isomorphic.

After this theorem, we can stop being so careful about our notation when dealing with these cyclic groups. Since they are isomorphic, any purely group theoretic question asked of them will give identical answers. We will use whichever version of the group lends itself to the question at hand, which is usually  $\mathbb{Z}_n$ .

Since an isomorphism between two groups means they are the same, we should check that the term behaves appropriately.

*Theorem 3.76.* Let  $G$ ,  $H$ , and  $K$  be groups. Then

1.  $G \cong G$ ;
2. if  $G \cong H$  and  $H \cong K$ , then  $G \cong K$ ; and
3. if  $G \cong H$ , then  $H \cong G$ .

An isomorphism of a group to itself just permutes the elements of the group while preserving the binary operation.

*Theorem 3.77.* Let  $G$  be a group with an element  $g$ . Define  $\phi_g : G \rightarrow G$  by  $\phi_g(h) = ghg^{-1}$ . Then  $\phi_g : G \rightarrow G$  is an isomorphism, called **conjugation by  $g$** .

One way to tell whether a homomorphism is an isomorphism is to look at its kernel and its image. The next theorem tells us that it is enough to check that  $e_H$  has only one preimage under  $\phi : G \rightarrow H$  to know that the whole function is injective.

*Theorem 3.78.* Let  $\phi : G \rightarrow H$  be a homomorphism. Then  $\phi$  is a monomorphism if and only if  $\text{Ker}(\phi) = \{e_G\}$ . In particular,  $\phi$  is an isomorphism if and only if  $\text{Im}(\phi) = H$  and  $\text{Ker}(\phi) = \{e_G\}$ .

The modular arithmetic groups give us many good examples of homomorphisms and isomorphisms.

*Theorem 3.79.* The map  $\phi : \mathbb{Z}_n \rightarrow \mathbb{Z}_n$  defined by  $\phi([a]_n) = [n - a]_n$  is an isomorphism.

*Theorem 3.80.* Let  $k$  and  $n$  be natural numbers. The map  $\phi : \mathbb{Z}_n \rightarrow \mathbb{Z}_n$  defined by  $\phi([a]_n) = [ka]_n$  is a homomorphism.

*Exercise 3.81.* Make and prove a conjecture that gives necessary and sufficient conditions on the natural numbers  $k$  and  $n$  to conclude that  $\phi : \mathbb{Z}_n \rightarrow \mathbb{Z}_n$  defined by  $\phi([a]_n) = [ka]_n$  is an isomorphism. Use this to show that there are several different isomorphisms  $\phi : \mathbb{Z}_{12} \rightarrow \mathbb{Z}_{12}$ .

*Exercise 3.82.* Use conjugation to find two different isomorphisms  $\phi : D_4 \rightarrow D_4$ . Why does conjugation not give any interesting isomorphisms from  $\mathbb{Z}_n$  to itself?

The isomorphisms above are all very straightforward, either from a group to itself or between groups obviously based on the same set (like  $C_n$  and  $\mathbb{Z}_n$ ). But isomorphisms need not be between a group and itself.

*Exercise 3.83.* Let  $H = \{[\text{just put the triangle back where you found it}], [\text{rotate counterclockwise by } 120^\circ], [\text{rotate counterclockwise by } 240^\circ]\}$ , which is a subgroup of  $D_3$ . Then  $H$  is isomorphic to the cyclic group  $C_3$ .

*Theorem 3.84.* For each  $n$ , the symmetries of a regular  $n$ -gon,  $D_n$ , has a subgroup isomorphic to  $C_n$ .

One strategy of mathematical exploration is to find the most general or the most comprehensive examples of a mathematical object and study it. The previous theorem says that the symmetry groups of the regular polygons contain the finite cyclic groups as subgroups. Similarly, the symmetric groups contain the symmetry groups of the regular polygons.

*Theorem 3.85.* For any  $n$ , the symmetric group  $S_n$  has subgroups isomorphic to  $D_n$  and  $C_n$ .

The previous theorem is a little misleading, since it makes it seem like only the groups corresponding to the same  $n$  have any relationships.

*Exercise 3.86.* Find 40 different subgroups of  $S_6$  isomorphic to  $C_3$ .

*Exercise 3.87.* For any  $m > n$ , find a monomorphism  $\phi : S_n \rightarrow S_m$ .

And the crowning theorem tells us that *every* group is a subgroup of a symmetric group.

*Theorem 3.88.* Let  $G$  be a group. Then for some set  $X$ , there is a subgroup  $H$  of  $\text{Sym}(X)$  such that  $G$  is isomorphic to  $H$ . If  $G$  is finite, then  $X$  can be chosen to be finite.

### 3.7 Sizes of Subgroups and Orders of Elements

*Definition.* Let  $H$  be a subgroup of a group  $G$ . Then the **left coset of  $H$  by  $g$**  is the set of all elements of the form  $gh$  for all  $h \in H$ . This left coset is written  $gH = \{gh \mid h \in H\}$ . Right cosets are defined similarly.

The notation in the previous definition works well when the binary operation  $*$  of the group  $G$  is written multiplicatively, like  $g * h = gh$ , for example in  $D_n$  or  $S_n$ . But when the binary operation is written additively, with a plus sign, this notation can be confusing. So when writing the cosets of an additive group we use a  $+$  notation for the cosets. For example, consider the group  $(\mathbb{Z}, +)$ . Then the cosets of the subgroup  $3\mathbb{Z}$  are written  $\{0 + 3\mathbb{Z}, 1 + 3\mathbb{Z}, 2 + 3\mathbb{Z}\}$ .

*Exercise 3.89.* Consider  $H = \{[\text{do nothing}], [\text{flip over the vertical line}], [\text{flip over the horizontal line}], [\text{rotate } 180^\circ \text{ counterclockwise}]\}$ , a subgroup of  $D_4$ , and  $K = \{[0]_{12}, [3]_{12}, [6]_{12}, [9]_{12}\}$ , a subgroup of  $\mathbb{Z}_{12}$ . Write out the cosets of  $H$  and  $K$  in their respective groups.

*Lemma 3.90.* Let  $H$  be a subgroup of  $G$  and let  $g$  and  $g'$  be elements of  $G$ . Then the cosets  $gH$  and  $g'H$  are either identical (the same subset of  $G$ ) or disjoint.

Recall that  $|G|$  denotes the order of the group  $G$ .

*Theorem 3.91 (Lagrange).* Let  $G$  be a finite group with subgroup  $H$ . Then  $|H|$  divides  $|G|$ .

Since the order of a subgroup divides the order of the group, it is natural to define a term that records how many times the order of the subgroup divides the order of the group.

*Definition.* Let  $H$  be a subgroup of a group  $G$ . Then the **index of  $H$  in  $G$**  is the number of distinct (left) cosets of  $H$ . We write this index as  $[G : H]$ .

*Theorem 3.92.* Let  $G$  be a finite group with a subgroup  $H$ . Then  $[G : H] = |G|/|H|$ . In particular, the number of left cosets of  $H$  is equal to the number of right cosets of  $H$ .

Lagrange's Theorem has many implications. One is that the order of each element must divide the order of the group.

*Corollary 3.93.* Let  $G$  be a finite group with an element  $g$ . Then  $o(g)$  divides  $|G|$ .

*Corollary 3.94.* If  $p$  is a prime, then  $C_p$  has no non-trivial subgroups.

### 3.8 Normal Subgroups

In general, a left coset  $gH$  may or may not be equal to the right coset  $Hg$ .

*Exercise 3.95.* Find a subgroup  $H$  of  $D_3$  and an element  $g$  of  $D_3$  such that  $gH$  is not equal to  $Hg$ .

Although in general a left coset  $gH$  may or may not be equal to the right coset  $Hg$ , when  $K$  is the kernel of a homomorphism,  $gK$  always equals  $Kg$ .

*Theorem 3.96.* Let  $G$  and  $H$  be groups and let  $\phi : G \rightarrow H$  be a homomorphism. Then for every element  $g$  of  $G$ ,  $gKer(\phi) = Ker(\phi)g$ .

It is useful to give a name to those subgroups, like kernels of homomorphisms, with the property that each right coset is equal to the corresponding left coset.

*Definition.* A subgroup  $K$  of a group  $G$  is **normal**, denoted  $K \triangleleft G$ , if and only if for every element  $g$  in  $G$ ,  $gK = Kg$  or, equivalently,  $K = gKg^{-1}$  ( $= \{gkg^{-1} | k \in K, \text{ and } g \in G\}$ ).

We can reformulate the previous theorem using this new vocabulary: if  $\phi : G \rightarrow H$  is a homomorphism, then  $Ker(\phi)$  is a normal subgroup of  $G$ .

In abelian groups, all subgroups are normal.

*Theorem 3.97.* Let  $K$  be a subgroup of an abelian group  $G$ . Then  $K$  is a normal subgroup.

An equivalent characterization of normal subgroups is often useful. Recall the definition of conjugation by  $g$  above. The reformulated characterization in the next theorem is often called " $K$  is closed under conjugation by every element  $g \in G$ ".

*Theorem 3.98.* A subgroup  $K$  of a group  $G$  is normal if and only if for every  $g \in G$ ,  $gKg^{-1} \subseteq K$ .

Recall that the center of a group  $G$ ,  $Z(G)$ , is the set of all the elements of  $G$  that commute with every element of  $G$ . The center of a group is always a normal subgroup.

*Theorem 3.99.* Let  $G$  be a group. Then  $Z(G) \triangleleft G$ .

### 3.9 Quotient Groups

We have discussed at least two methods for creating new groups from we already have: one source was to start with a group and locate its subgroups, another was to start with two groups and take their product. A third method for creating new groups from old involves quotients, which will be associated with normal subgroups.

If  $\phi : G \rightarrow H$  is a homomorphism, then we already know that  $\text{Im}(\phi)$  is a group. It turns out that we can think of the elements of this group as the cosets of the kernel of  $\phi$ . This will allow us to create a group out of  $G$  and a normal subgroup  $K$  without mentioning a homomorphism.

*Exercise 3.100.* Let  $\phi : G \rightarrow H$  be an epimorphism. Let  $h \in H$  and show that  $\phi^{-1}(h) = g\text{Ker}(\phi)$  for some  $g \in G$ .

Of course, we already knew that  $\text{Im}(\phi)$  was a group. However, there is an intrinsic way to create a group whose elements are the cosets of a normal subgroup.

*Theorem 3.101.* Let  $K$  be a normal subgroup of a group  $G$ , and let  $G/K$  be the left cosets of  $K$  in  $G$ . Define the binary operation  $\hat{*}$  on  $G/K$  by  $gK \hat{*} g'K = (gg')K$ . Then  $(G/K, \hat{*})$  is a group, and  $|G/K| = [G : K]$ .

When  $K$  is a normal subgroup of  $G$ , the group  $G/K$  described above is called a **quotient group** and read “ $G \bmod K$ ”.

*Exercise 3.102.* Explain the necessity of the normality hypothesis in the definition of quotient groups. Give an example of a group  $G$  with a subgroup  $H$  such that the cosets of  $H$  do not form a group using the operation defined in Theorem 3.101.

Since there are several different binary operations floating around in a quotient group, the notation can get a little confusing, so let’s look at two examples.

- (1) Let  $H$  be a normal subgroup of  $D_n$  and let  $s_1$  and  $s_2$  be elements of  $D_n$ . Then  $s_1H \hat{\circ} s_2H = (s_1 \circ s_2)H$ .
- (2) The cosets of the normal subgroup  $3\mathbb{Z}$  of  $\mathbb{Z}$  are written  $\{0 + 3\mathbb{Z}, 1 + 3\mathbb{Z}, 2 + 3\mathbb{Z}\}$ . And in  $\mathbb{Z}/3\mathbb{Z}$ ,  $(1 + 3\mathbb{Z}) \hat{+} (1 + 3\mathbb{Z}) = (1 + 1) + 3\mathbb{Z} = 2 + 3\mathbb{Z}$ .

*Theorem 3.103.* For any natural number  $n$ ,  $n\mathbb{Z}$  is a normal subgroup of  $\mathbb{Z}$ . Furthermore,  $\mathbb{Z}/n\mathbb{Z} \cong C_n \cong \mathbb{Z}_n$ .

Taking a quotient of an object (in this case a group) by a sub-object (in this case a normal subgroup) is like identifying the elements of the sub-object with each other. In the case of a quotient group, we form a new

group where the elements in each coset become one element. This technique is rich and is repeated often in mathematics.

To get accustomed to quotient groups, let's look at the quotient groups that arise from  $S_4$ .

*Exercise 3.104.* Find all normal subgroups  $K$  of  $S_4$ . For each such  $K$  show that  $S_4/K$  is isomorphic to  $\{e\}$ ,  $C_n$  for some  $n$ ,  $D_n$  for some  $n$ , or  $S_n$  for some  $n$ .

Sometimes the structure of the quotient group can give us information about the whole group.

*Theorem 3.105.* Suppose that  $G/Z(G)$  is a cyclic group. Then  $G$  is abelian. (So  $Z(G) = G$ .)

Conversely, sometimes knowledge about the normal subgroup gives us information about the character of the quotient group.

*Theorem 3.106.* Let  $K$  be a normal subgroup of a group  $G$ . Then  $G/K$  is abelian if and only if  $K$  contains  $\{ghg^{-1}h^{-1} | g, h \in G\}$ .

Because of its role in making quotient groups commute, the subgroup of  $G$  generated by elements of the form  $ghg^{-1}h^{-1}$ , that is

$$[G, G] = \langle \{ghg^{-1}h^{-1} | g, h \in G\} \rangle,$$

is called the **commutator subgroup** of  $G$ .

*Theorem 3.107.* Let  $G$  be a group, then  $[G, G]$  is a normal subgroup of  $G$ . Thus  $G/[G, G]$  is an abelian group, called the **abelianization** of  $G$ .

Quotient groups naturally arose from investigating images of homomorphisms. It turns out that every quotient group is the image of a homomorphism.

*Theorem 3.108.* Let  $G$  be a group and  $K$  be a normal subgroup of  $G$ . Then there is an epimorphism  $q : G \rightarrow G/K$  with  $\text{Ker}(q) = K$ . So a subgroup is normal if and only if it is the kernel of a homomorphism.

One of the most useful and important theorems in group theory relates the image of a homomorphism with a quotient group.

*Theorem 3.109 (First Isomorphism Theorem).* For any homomorphism,  $\phi : G \rightarrow H$ , the image of  $\phi$  is isomorphic to the quotient group  $G \text{ mod } \text{Ker}(\phi)$ , or using the notation for isomorphism:  $\text{Im}(\phi) \cong G/\text{Ker}(\phi)$ .

*Corollary 3.110.* For any epimorphism,  $\phi : G \rightarrow H$ ,  $H \cong G/\text{Ker}(\phi)$ .

The First Isomorphism Theorem allows us to determine the structure of many groups and their subgroups. There are other Isomorphism Theorems, but they are less fun and more technical.

*Theorem 3.111.* For each natural number,  $n$ , there is a unique cyclic group with order  $n$ .

*Theorem 3.112.* Let  $m$  and  $n$  be relatively prime integers. Then  $\mathbb{Z}_m \times \mathbb{Z}_n \cong \mathbb{Z}_{mn}$ .

One of the central theorems in the study of abelian groups relates them to products of cyclic groups. The proof of the following theorem is rather involved, but can be done by finding an epimorphism from  $\mathbb{Z}^n$  ( $= \mathbb{Z} \times \cdots \times \mathbb{Z}$   $n$  times) to your group and cleverly describing the kernel.

*Lemma 3.113.* Let  $G$  be a finitely generated abelian group, generated by  $n$  elements  $\{g_1, g_2, \dots, g_n\}$ . Then  $\phi : \mathbb{Z}^n \rightarrow G$  defined by  $(a_1, a_2, \dots, a_n) \mapsto g_1^{a_1} g_2^{a_2} \cdots g_n^{a_n}$  is an epimorphism.

If we understood the kernel of the epimorphism from the previous lemma, we could use the First Isomorphism Theorem to prove the next theorem, called the Fundamental Theorem of Finitely Generated Abelian Groups. The following is a difficult theorem.

*Theorem\** 3.114. Every finitely generated abelian group is isomorphic to a direct product of cyclic groups.

This theorem is considered fundamental because it describes the structure of every finitely generated abelian group as being built up as a product of cyclic groups. As we mentioned above, decomposing a group into products tells us a lot about it.

### 3.10 More Examples

Before we conclude our tour of group theory, let's describe a few additional groups. Like the symmetric groups, the elements of the following example are functions.

*Exercise 3.115.* Let  $M = \{f(x) = \frac{ax+b}{cx+d} \mid a, b, c, d \in \mathbb{R}, ad - bc = 1\}$ .

1. Show that if  $f(x)$  and  $g(x)$  are functions in  $M$ , then  $f \circ g(x)$ , their composition, is also of that form (i.e., in  $M$ ).
2. Given  $f(x) = \frac{ax+b}{cx+d}$ , find another function in  $M$ ,  $h(x)$ , such that  $f \circ h(x) = h \circ f(x) = x = \frac{1x+0}{0x+1}$ .



*Theorem 3.116.* The set  $M = \{f(x) = \frac{ax+b}{cx+d} | a, b, c, d \in \mathbb{R}, ad - bc = 1\}$  with the binary operation of composition forms a group.

Here is a group whose elements are matrices.

*Theorem 3.117.* The set of  $2 \times 2$  matrices with real number entries and determinant equal to 1, written  $SL_2(\mathbb{R})$ , along with the operation of matrix multiplication is a group.

*Theorem 3.118.* The groups  $M$  and  $SL_2(\mathbb{R}) / \langle \begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix} \rangle$  are isomorphic.

*Exercise 3.119.* Realize  $D_4$  as a subgroup of  $SL_2(\mathbb{R})$  by writing each transformation as a matrix. In other words, find an injective homomorphism  $\phi : D_4 \rightarrow SL_2(\mathbb{R})$ , which is the same as finding a subgroup of  $SL_2(\mathbb{R})$  isomorphic to  $D_4$ .

The self-isomorphisms of a group form a group of their own.

*Definition.* Let  $G$  be a group. Then  $Aut(G)$  is the set of isomorphisms from  $G$  to itself, which comes with composition of functions as a binary operation. The elements of  $Aut(G)$  are called **automorphisms** of  $G$ .

*Theorem 3.120.* Let  $G$  be a group. Then  $(Aut(G), \circ)$  is a subgroup of  $Sym(G)$ . In particular,  $Aut(G)$  is a group.

*Theorem 3.121.* If  $p$  is a prime number, then  $Aut(\mathbb{Z}_p) \cong \mathbb{Z}_{p-1}$ .

We've actually already seen some automorphisms. Let  $g$  be an element of the group  $G$ . Then we defined  $\phi_g : G \rightarrow G$ , conjugation by  $g$ , by  $\phi_g(h) = ghg^{-1}$  for every  $h \in G$ .

*Theorem 3.122.* Let  $G$  be a group. For each element  $g \in G$ ,  $\phi_g$  is an automorphism of  $G$ . Furthermore,  $\Phi : G \rightarrow Aut(G)$  defined by  $\Phi(g) = \phi_g$  is a homomorphism.

The preceding theorem provides us with a bunch of automorphisms for non-abelian groups, but for abelian groups conjugation is just the trivial automorphism. The next theorem tells us that every group with at least three elements has at least one non-trivial automorphism.

*Theorem\** 3.123. Let  $G$  be a group with more than 2 elements. Then  $|Aut(G)| > 1$ .

### 3.11 Groups in Action

Although we introduced groups as an abstract structure, they actually appear in many different applications. In fact, historically, groups grew out of observations about collections of permutations satisfying the Cancellation

Law. These groups that “do things” can be thought of as *acting on a set*. We’ve already seen a few groups that act on sets in that sense, for example,  $D_n$  transforms the regular  $n$ -gon, and  $Sym(X)$  permutes the elements of  $X$ . Since  $Sym(X)$  includes *any* permutation of the set  $X$ , it is our most general example of a group acting on a set. When thinking about a group “doing something”, we want the elements of the group to be associated with permutations, but in a way that respects the group structure. We’ll first give a formal definition of this idea and then explain what it means.

*Definition.* Let  $G$  be a group. Then an **action** of  $G$  on a set  $X$  is a homomorphism  $\phi : G \rightarrow Sym(X)$ . We say that  $G$  *acts on*  $X$  by  $\phi$ .

At first, the term “action” might seem a little confusing. However, a map  $\phi$  from  $G$  to  $Sym(X)$  allows us to think of each element  $g$  of  $G$  as a permutation, namely the permutation  $\phi(g) \in Sym(X)$ . With that relationship,  $G$  is a collection of permutations, each of which acts on the set in the usual sense of being a permutation. And in this form, we can use our intuition about permuting objects to answer group-theoretic questions.

We have already seen several examples of actions. In Theorem 3.88, we used the fact that a group acts on itself by left multiplication to see that any group can be realized as a subgroup of  $Sym(G)$ .

*Theorem 3.124.* Let  $G$  be a group. For each  $g \in G$ , define  $\lambda_g : G \rightarrow G$  by  $\lambda_g(h) = gh$ . Then  $\Lambda : G \rightarrow Sym(G)$  defined by  $g \mapsto \lambda_g$  is an action of  $G$  on  $G$ .

Similarly,  $G$  acts on the cosets of a subgroup  $H$ .

*Theorem 3.125.* Let  $H$  be a subgroup of a group  $G$  and let  $L = \{gH | g \in G\}$  be the set of left cosets of  $H$ . Then the function  $\phi : G \rightarrow Sym(L)$  defined by  $\phi(g)(g'H) = (gg')H$  is an action of  $G$  on  $L$ .

We have already seen a second example of a group  $G$  acting on itself in Theorem 3.122 when we defined the homomorphism  $\Phi : G \rightarrow Aut(G) \subset Sym(G)$  that associated each element  $g$  of  $G$  with the automorphism: conjugate by  $g$ . In other words,  $G$  acts on itself by conjugation.

When we begin to explore the idea of group actions, two ideas arise about how the elements of  $G$  are moving the elements of  $X$  around. The first natural question is, “Where does an element  $x \in X$  go to under the permutations of  $G$ ?” The second natural question is, “For each element of  $X$ , what elements of  $G$  leave it fixed?” These questions lead to two definitions.

*Definitions.* 1. Let  $G$  be a group with an action  $\phi : G \rightarrow Sym(X)$  and let  $x \in X$ . The set of group elements that fix  $x$ , called the **stabilizer**

of  $x$ , is  $\text{Stab}(x) = \{g \in G \mid \phi(g)(x) = x\}$ .

2. The **orbit** of  $x$  is  $\text{Orb}(x) = \{y \in X \mid y = \phi(g)(x) \text{ for some } g \in G\}$ , which is just the collection of elements that  $x$  gets mapped to by the action.

*Exercise 3.126.* Pick a non-trivial subgroup  $H$  of  $D_4$ , and consider  $G$  acting on the left cosets of  $H$ . For each coset  $gH$ , find its stabilizer and orbit under this action.

The reason for requiring that an action of  $G$  on  $X$  be a homomorphism from  $G$  to  $\text{Sym}(X)$ , instead of just any old function, is that we would like the action of an element  $g$  followed by the action of an element  $h$  to be the same as the action of the element  $hg$ . If we write out this condition, it is precisely what is required for the action to be a homomorphism. This condition guarantees that stabilizers are subgroups. In other words, an action respects the group's structure.

*Theorem 3.127.* Let  $G$  be a group and  $\phi : G \rightarrow \text{Sym}(X)$  an action of  $G$  on  $X$ . If  $x \in X$ , then  $\text{Stab}(x)$  is a subgroup of  $G$ .

One of the neatest features about group actions is that there is a basic relationship between the number of places an element of  $X$  goes to under the action of  $G$  and the number of elements of  $G$  that leave it fixed.

*Theorem 3.128.* Let  $G$  be a finite group acting on a set  $X$ . Then for any  $x \in X$ ,

$$|\text{Stab}(x)| \cdot |\text{Orb}(x)| = |G|.$$

The orbits partition  $X$ .

*Lemma 3.129.* Let  $G$  be a group acting on a set  $X$ . Then for two elements  $x, y \in X$ , either  $\text{Orb}(x) = \text{Orb}(y)$  or they are disjoint.

*Theorem 3.130.* Let  $G$  be a finite group acting on a finite set  $X$ . Then  $|X| = \sum \frac{|G|}{|\text{Stab}(x_i)|}$ , where the sum is taken over one element  $x_i$  from each distinct orbit.

These theorems that relate the sizes of the group, the set, the stabilizers, and the orbits give clever methods for gaining insights into the structure of finite groups. For example, Cauchy proved that if a prime  $p$  divides the order of a group, then the group has an element of order  $p$ . Let's prove it using a group action.

*Lemma 3.131.* Let  $G$  be a finite group. Let  $S(n)$  be the set of all  $n$ -tuples of elements of  $G$  such that the product of those  $n$  elements in order equals the identity element. That is,  $S(n) = \{(g_1, g_2, g_3, \dots, g_n) \mid g_i \in G \text{ and } g_1 g_2 g_3 \dots g_n = e_G\}$ . Then  $|S(n)| = |G|^{n-1}$ .

*Lemma 3.132.* Let  $G$  be a finite group,  $p$  be a prime that divides the order of  $G$ , and  $S(p) = \{(g_1, g_2, g_3, \dots, g_p) | g_i \in G \text{ and } g_1 g_2 g_3 \dots g_p = e_G\}$ . Let  $\phi : \mathbb{Z}_p \rightarrow \text{Sym}(S(p))$  be defined by cyclic permutation, that is,

$$\phi([i]_p)((g_1, g_2, g_3, \dots, g_p)) = (g_{1+i}, g_{2+i}, \dots, g_{p+i})$$

where the subscripts are interpreted mod  $p$ . Then  $\text{Stab}((e_G, e_G, e_G, \dots, e_G)) = \mathbb{Z}_p$ .

*Theorem 3.133 (Cauchy).* Let  $G$  be a finite group and  $p$  be a prime that divides the order of  $G$ , then  $G$  has an element of order  $p$ .

Recall that Lagrange's Theorem stated that the order of any subgroup of a finite group divided the order of the group. There is a partial converse of this theorem that can be proved using the ideas of groups acting on sets. Cauchy's Theorem is a special case of this more general theorem.

*Theorem\* 3.134 (Sylow).* Let  $G$  be a finite group,  $p$  be a prime. If  $p^i$  divides the order of  $G$ , then  $G$  has a subgroup of order  $p^i$ .

One of the most fruitful sets on which a group can act is itself. We have already seen the action of conjugation, and we will now look at some further consequences of that action. This action is so important that its orbits and stabilizers have special names.

*Definitions.* 1. Let  $G$  be a group. For any element  $g \in G$ , define  $C_G(g) = \{h \in G | hgh^{-1} = g\}$ , called the **centralizer** of  $g$ .

2. Let  $g$  be an element of  $G$ ; then the **conjugacy class** of  $g$  is the set  $\{hgh^{-1} | h \in G\}$ .

*Exercise 3.135.* Describe the conjugacy classes in the symmetric groups,  $S_n$ .

*Corollary 3.136.* Let  $G$  be a finite group. Then  $|G| = \sum [G : C_G(g_i)]$ , where the sum is taken over one element,  $g_i$ , from each conjugacy class in  $G$ .

This last corollary can be used to show that certain groups have non-trivial centers, without producing specific elements that commute!

*Theorem 3.137.* Let  $G$  be a group such that  $|G| = p^n$  for some prime  $p$ . Then  $|Z(G)| \geq p$ .

## 3.12 The Man Behind the Curtain

Many people mistakenly believe that mathematics is arbitrary and magical, or at least that there is some secret knowledge that math teachers have but

won't share with their students. Mathematics is no more magical than the Great and Powerful Wizard of Oz, who was just a man behind a curtain. The development of mathematics is directed by a few simple principles and a strong sense of aesthetics. To develop the ideas of graph theory and group theory we followed a path of guided discovery. Let's look back on the journey and let the guiding strategies emerge from behind the curtain.

We started with examples; graphs and groups did not appear fully formed. Those ideas emerged from pinning down the essential features and commonalities of specific examples. We distilled those essentials into definitions. Definitions focus our attention on some features of our generative examples, but other choices for emphasis are possible, and making other choices would lead to other mathematics. For example, focusing on other additive and multiplicative properties of the reals or rationals leads to the definitions of other algebraic structures besides groups, including objects called rings and fields. Similarly, when defining the concept of a graph, if we were interested in questions involving directed connections, we would be led to the subject of directed graphs.

After isolating the concepts of graphs and groups, we explored the implications of our definitions. We created concepts concerning graphs and groups that allowed us to differentiate special subtypes of graphs and groups and to find and express theorems about their structure. Our exploration involved defining sub-objects (like subgroups), isolating the meaning of sameness (like isomorphisms), and developing a concept of functions that preserve the structure (like homomorphisms).

In making decisions about what definitions are appropriate and what theorem statements are valuable, the aesthetics of mathematics plays a significant role. Ideally a definition should capture and clarify a concept and a theorem should illuminate a relationship so that we get a satisfying sense of insight. A theorem should be as clean and general as possible.

We will see these strategies for creating mathematics and this sense of mathematical aesthetics repeated and refined in our exploration of other abstract mathematics in the chapters ahead.



## Chapter 4

# Calculus

### 4.1 Know your Whereabouts

Well-known prankster Zeno was notorious for creating elaborate paintings in prominent places throughout the city. As far as the city officials were concerned, this graffiti was vandalization of public property, so the police set up a camera to take pictures of the his favorite wall to catch him in the act. But the camera broke at the moment that Zeno was going to touch paint to concrete, and the case was ruined. Fortunately, two cops, officer Isaac and officer Gottfried, stepped forward to take over the investigation. It shocked their boss when, rather than trying to find more evidence, the officers claimed they had all the photos they needed and went to arrest our hapless “hero”, Zeno. Here’s what they said when they picked him up:

Officer Isaac: Zeno, at 12:59am your brush was 100ft from the wall, as seen in this photo.

Office Gottfreid: At 12:59.5am your brush was 25ft from the wall, as seen in this second photo.

Office Isaac: At 12:59.75am your brush was  $\frac{25}{4}$ ft from the wall, as seen in this third photo.

Officer Gottfreid: In fact, at  $(\frac{1}{2})^n$  minutes before 1am, your brush was  $\frac{100}{4^n}$ ft from the wall. This is an infinite number of claims, but we had a special camera that could take this infinite sequence of pictures in the moments leading up to 1am. (It’s no surprise that this camera overheated and stopped working right when it did, huh?)

Officer Isaac: Although we do not have a photograph at 1am showing you actually painting, these photos show that your brush did make contact

with the wall at 1am because the positions of the brush at times arbitrarily close to 1am converge to touching the wall at 1am.

Zeno: Converge? What does *converge* mean?

Officer Gottfried: Just think it over for 150 years and you'll understand. You'll have plenty of time while you're in jail.

Zeno: Maybe you'll be the ones with too much time on your hands when they revoke your badges for this stunt.

The moral of this story is that the motion of an object in space is predictable; objects do not jump or teleport when moving. If you know Zeno's position except at one instant in time, you know where he is at all times.

Suppose you know the position of a certain particle at all times right before and right after time  $t_0$ , but not right at  $t_0$ . You can think of this as a movie that is missing a single frame in the middle. There is only one way to insert the missing frame so that the particle's motion appears smooth. To find out where to put the particle in the missing frame, you can do exactly what Isaac and Gottfried did to find Zeno's position at time  $t_0$ , namely pick a point that you think is correct and make sure that the positions on the nearby frames are getting arbitrarily close to that point.

This process is a little complicated, so let's abstract and simplify a little. Notice that, although Zeno was moving in a 3-dimensional world, we could use a single number, his distance to the wall, to represent his position. So for each small time interval before 1pm we also have a number representing his position. If you make a list out of these numbers for his positions at 1 minute before 1pm,  $\frac{1}{2}$  a minute before 1pm,  $\frac{1}{4}$  of a minute before 1pm, and so forth, then these numbers are 'approaching' a particular number, which we will call  $\ell$ . Then we are asserting that Zeno must be at position  $\ell$  at 1pm.

$$(100, 25, 6.25, 1.56, \dots) = (100, \frac{100}{4}, \frac{100}{16}, \frac{100}{64}, \dots) = (\frac{100}{4^{n-1}} | n \in \mathbb{N}) \rightarrow \ell = 0$$

It will take us quite a bit of work to turn this notion of 'approaching' into a precise definition of *convergence*.

## 4.2 Convergence

Pinning down the idea of convergence required mathematicians more than 150 years. The challenge is to describe what it means for an infinite *sequence*



of numbers to *converge* to a single number, called the *limit*. As with many great mathematical insights, the solution came from considering an easier problem. First, we will glean some common properties of sequences that, at least intuitively, ‘approach’ a fixed number; and then, rather than trying to pin down *exactly* what number an arbitrary sequence ‘approaches’, we start by trying to do better than a certain error. Finally, we investigate an appropriate way to say that we can do better than any given error. Of course, we first need a precise definition of the objects we’ll be studying.

*Definition.* A **sequence** is an ordered list of real numbers indexed by the natural numbers  $(a_1, a_2, a_3, \dots) = (a_n | n \in \mathbb{N}) = (a_n)_{n \in \mathbb{N}}$ .

In terms of this new vocabulary, we are interested in defining and understanding when a sequence ‘approaches’ a fixed number  $\ell$ . Instead of trying to define exactly what this means, let’s start by observing some things that had better be true about any definition that captures the notion of ‘approaching’.

Observation 1: If the sequence  $S = (a_n | n \in \mathbb{N})$  ‘approaches’ a number  $\ell$ , then ‘eventually’ the terms of  $S$  had better be ‘very close’ to the number  $\ell$ .

There are two major parts to this observation that we must investigate. What is the precise meaning of ‘eventually’, and what is the best definition of ‘very close’? We start by trying to make the notion of ‘very close’ more precise, but in the end, the precise definitions of these two ideas will depend on each other. In particular, the observation highlights our need to quantify how far two real numbers are from each other, so we start there.

*Definition.* Let  $x$  and  $y$  be two real numbers. Then the **distance** from  $x$  to  $y$  is defined as  $\|y - x\|$ , the absolute value of  $y - x$ . If  $\|y - x\| < \varepsilon$ , then we say that  $y$  is **within a distance of  $\varepsilon$  from  $x$** .

Recall the definition of the absolute value function: if  $a$  is a non-negative real number then  $\|a\| = a$ , and if  $a$  is a negative real number then  $\|a\| = -a$ . It follows that, for any real number  $a$ ,  $0 \leq \|a\|$ . Furthermore,  $\|a\| = 0$  if and only if  $a = 0$ . We will be using the absolute value function in almost every proof in this chapter, so we should warm up with a few basic properties. As always, carefully use the definition of the absolute value function, not any preconceived notion of how it works.

*Lemma 4.1.* Let  $a, b \in \mathbb{R}$ .

1.  $\|ab\| = \|a\|\|b\|$

$$2. a \leq \|a\|$$

$$3. \|-a\| = \|a\|$$

In particular, the last equality in Lemma 4.1 tells us something about our definition of distance: for any two real numbers,  $x$  and  $y$ ,  $\|x - y\| = \|y - x\|$ , so we may call this quantity the distance between  $x$  and  $y$  (instead of the distance from  $x$  to  $y$ ).

For this to be a good notion of distance, it should satisfy one additional property, called the *Triangle Inequality*. Essentially, this inequality captures the notion that it is always shorter to go directly from point  $A$  to point  $B$  than it is to go first from  $A$  to  $C$  then from  $C$  to  $B$ .

*Lemma 4.2* (Triangle Inequality). Let  $x, y, z$  be real numbers. Then  $\|y - x\| \leq \|y - z\| + \|z - x\|$ . (Hint: Note that  $\|y - x\|^2 = (y - x)^2$ .)

The triangle inequality has many equivalent formulations, and they are also useful.

*Corollary 4.3.* Let  $a, b \in \mathbb{R}$ . Then  $\|a\| - \|b\| \leq \|a + b\| \leq \|a\| + \|b\|$ .

We now have the language to talk about the distance between two numbers, but what does it mean for two points to be ‘very close’ to each other? The answer is subtle because ‘very close’ is a relative term. The distance between my home and the grocery store is pretty small; I can drive there in under 3 minutes. But to an ant, the distance is enormous; an ant could probably walk for days before reaching the grocery store from my house. Moreover, when my car is in the shop, I have to walk to the grocery store, and carrying groceries in the heat that far makes the distance seem insurmountably large. The point is, for any two distinct points, there is a perspective in which they appear quite far apart.

So, to say that two points are ‘very close’, we must first choose a perspective, which means that we must set an allowable threshold for points to be ‘very close’. If we decide that two points are ‘very close’ if a person can drive between them in under 15 minutes, then my house and the grocery store are ‘very close’. However, if we decide that two points are ‘very close’ only if an ant can walk between them in less than a day, then my house and the grocery store are not ‘very close’. Similarly, if we say that two real numbers within a distance of 0.5 from each other are ‘very close’ to each other, then 2 and 2.1 are ‘very close’ to each other. But if we have stricter standards and require the numbers to be within a distance of 0.001 from each other, then 2 and 2.1 are not ‘very close’ to each other.

Sadly, permanently fixing a certain distance to be the cut-off for ‘very close’ does not produce a reasonable notion of ‘approaching’, as you will show in the following exercise.

*Exercise 4.4.* Suppose that we decreed two real numbers to be ‘very close’ to each other if the distance between them is less than 0.1. Describe a specific sequence  $S = (a_n | n \in \mathbb{N})$  whose terms are all ‘very close’ to 5 but whose terms do not ‘approach’ 5. Of course you must also include a description of your notion of ‘approaching’ in your explanation, or at least explain what property your sequence has that is contrary to the intuitive notion of ‘approaches’ 5.

So we see that we cannot fix a distance as ‘very close’ before defining the notion of ‘approaching’. Instead, for a sequence  $S$  to ‘approach’  $\ell$ , it’s terms must be ‘very close’ to  $\ell$ , regardless of which notion of ‘very close’ is chosen. So the next logical thing to try is requiring the sequence’s terms to be ‘very close’ for *all* choices of ‘very close’?

*Exercise 4.5.* Let  $S = (a_n | n \in \mathbb{N})$  be a sequence. Suppose all terms in  $S$  are ‘very close’ to the number  $\ell$ , regardless of which distance is chosen as the cut-off for ‘very close’. Describe  $S$  fully by finding out exactly the value of each  $a_n$ .

The previous exercise shows that requiring every term to be ‘very close’ for every notion of ‘very close’ is far too restrictive. So let’s return to the sequence of positions from Zeno’s graffiti escapades for inspiration:

$$S = (100, \frac{100}{4}, \frac{100}{16}, \dots) = (\frac{100}{4^{n-1}} | n \in \mathbb{N}).$$

We began this search to try to formalize the intuitive notion that this sequence ‘approaches’ 0, which corresponds to the distance between brush and wall being decreasing to 0 (which, in turn, corresponds to the pain meeting the wall). We have not captured even our motivating example, so we must keep looking to find an appropriate notion. If we require that all terms be ‘very close’ to 0 for any definition of ‘very close’, then we are in trouble, since any choice of ‘very close’ less than 100 causes a problem with the first term ( $\|100 - 0\| = 100$ ). But this was not a problem for our intuition. In terms of Observation 1, this restrictive requirement ignores the idea that the terms are ‘eventually’ close and instead requires them to be close from the outset. By trying to include the notion of ‘eventually’ we will be able to find a better notion of ‘approaching’.

Similarly, if  $S$  ‘approaches’ 0, then so should the sequence

$$S' = (17, 15, 3, 100, \frac{100}{4}, \frac{100}{16}, \dots).$$

The intuitive notion of ‘approaching’ has to do with only a *tail* of a sequence and can ignore any finite number of terms at the beginning. This example also foreshadows the idea that the sequence’s terms do not need to constantly decrease for our intuition to feel that it ‘approaches’ something.

*Definition.* Let  $S = (a_1, a_2, a_3, \dots)$  be a sequence. Then a **tail** of the sequence  $S$  is a set of the form  $\{a_m | m \geq M\}$ , where  $M$  is some fixed natural number. Note that  $S$  has a different tail for each choice of  $M$ , so it does not make sense to talk about the tail.

The key idea is that every property that must ‘eventually’ be true of a sequence can be stated in terms of tails of that sequence. In particular, if a sequence  $S$  ‘approaches’ a number  $\ell$ , then for any choice of ‘very close’ there is a tail of the sequence  $S$  such that all the terms in that tail are ‘very close’ to  $\ell$ .

Using the new notions of distance and tails, we can reformulate our observation more precisely.

Observation 1': If a sequence  $S = (a_n | n \in \mathbb{N})$  ‘approaches’ a number  $\ell$ , then for any cut-off for ‘very close’, there is a tail of  $S$  such that all terms in the tail are ‘very close’ to  $\ell$ . In other words, for each  $\varepsilon > 0$ , there should exist an  $M_\varepsilon \in \mathbb{N}$  such that for any  $k \geq M_\varepsilon$ ,  $\|a_k - \ell\| < \varepsilon$ .

This is a very technical idea, so let’s do a few computational exercises to get familiar with it.

*Exercise 4.6.* For each sequence below, you will be given a positive real number  $\varepsilon$  representing the cut-off for ‘very close’. Find a natural number  $M_\varepsilon$  such that all the terms beyond the  $M_\varepsilon^{\text{th}}$  term lie within the prescribed distance,  $\varepsilon$ , from  $\ell$ . It is not necessary to find the smallest value for  $M_\varepsilon$  for which this condition is true. As always, justify your answers.

1. Consider the sequence  $(1, \frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \frac{1}{5}, \dots) = (a_n = \frac{1}{n} | n \in \mathbb{N}) = (\frac{1}{n})_{n \in \mathbb{N}}$ . Find a natural number  $M_{0.03}$  such that, for every natural number  $k \geq M_{0.03}$ , each term  $a_k$  lies within a distance of  $\varepsilon = 0.03$  from  $\ell = 0$ .
2. Consider the sequence  $(1 - e^{-n} | n \in \mathbb{N})$ . Find a natural number  $M_{0.001}$  such that, for every natural number  $k \geq M_{0.001}$ , each term  $a_k = 1 - e^{-k}$  lies within a distance of  $\varepsilon = 0.001$  from  $\ell = 1$ .
3. Consider the sequence  $(\frac{(-1)^n}{n^2} | n \in \mathbb{N})$ . Find a natural number  $M_{0.0001}$  such that, for every natural number  $k \geq M_{0.0001}$ , each term  $a_k = \frac{(-1)^k}{k^2}$  lies within a distance  $\varepsilon = 0.0001$  from  $\ell = 0$ .

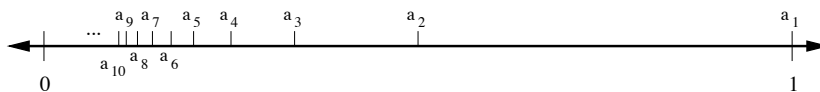
We have found a nice way of stating the property that “eventually the terms get very close to  $\ell$ ”, but what if our statement is too lenient, meaning too many sequences satisfy it? We should check that sequences that our intuition dictates do not approach a fixed number do not satisfy this condition, which involves articulating the negation of our statement.

*Exercise 4.7.* We have essentially defined the phrase “the sequence  $S$  ‘approaches’  $\ell$ ” as “for every cut-off for ‘very close’, there is a tail of  $S$  whose terms are all ‘very close’ to  $\ell$ ”. Using this provisional definition, write out a precise meaning for “the sequence  $S$  does not ‘approach’  $\ell$ ”.

In the previous exercise, you have articulated what it means for a sequence not to ‘approach’ a particular number  $\ell$ . But in practice you will be more interested in saying that a sequence doesn’t ‘approach’ *any* number. The hard part about checking this stronger condition is that you now need to show that it doesn’t work for *any* choice of  $\ell$ . This means, for each real number  $\ell$  you need to check that “the sequence  $S$  does not ‘approach’  $\ell$ ”. Of course, there are an infinite number of real numbers, so doing this for each value of  $\ell$  separately is not a viable option. Sometimes the proof can be done all at once, and sometimes you will need to break the argument into cases, depending on which value of  $\ell$  you are considering.

Up to this point, most of our sequences have been written as lists generated by formulas. Some people have strong intuition about such algebraic objects. Others have much more visual and geometric intuition, so we should find a good way to draw a sequence so that we may use this intuition. The terms in a sequence are real numbers, and we already have a good way to draw the real numbers: a line. So let’s represent the sequence on the real line.

For example, consider the sequence ( $a_n = \frac{1}{n} | n \in \mathbb{N}$ ). All of the terms in this sequence are between 0 and 1, so we should make sure to draw that part of the line very large. Then put a mark for each term in your sequence, and label it as follows.



Of course, you will not be able to draw the infinite number of terms in your sequence, just a representative sample of the first several terms. Knowing how many terms it takes to capture what the sequence is doing is a matter of experience.

Since we are trying to use the pictures to decide if the sequence is ‘ap-

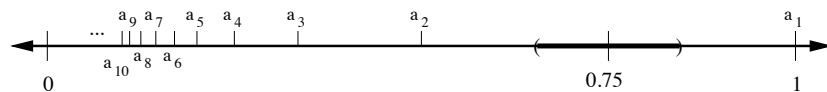
proaching' a particular number  $\ell$ , we should include  $\ell$  in the drawing as well. Also, given a distance  $\varepsilon$  that is the cut-off for 'very close', we would like to be able to draw the set of points that are very close to  $\ell$ . Fortunately, the set of numbers that are within  $\varepsilon$  from  $\ell$  is a very easy set to draw, called an *open interval*.

**Definition.** An **open interval** is a subset of the real numbers of the form  $\{r \in \mathbb{R} | a < r < b\}$  for two real numbers  $a < b$  and is called the "(open) interval from  $a$  to  $b$ " or the "(open) interval with endpoints  $a$  and  $b$ ". We will usually write  $(a, b)$  for the (open) interval from  $a$  to  $b$ , and call  $a$  and  $b$  the **endpoints** of the interval.

Note that the symbol " $(a, b)$ " could be an interval or an ordered pair, but the context should make it clear which is intended. And, as we were claiming above, the set of points with a distance of  $\varepsilon$  from  $\ell$  is an open interval:

$$\{x \in \mathbb{R} | |x - \ell| < \varepsilon\} = \{x \in \mathbb{R} | \ell - \varepsilon < x < \ell + \varepsilon\} = (\ell - \varepsilon, \ell + \varepsilon).$$

In other words, the set of numbers that are within a distance of  $\varepsilon$  from  $\ell$  is a line segment centered at  $\ell$ , not including the segments's endpoints. Usually, we draw this by putting parantheses on the number line at the endpoints and sometimes by shading the segment. For example, if we wanted to see if the sequence  $(a_n = \frac{1}{n} | n \in \mathbb{N})$  above is 'eventually' within a distance 0.1 from  $\ell = 0.75$ , we would add the following to the drawing.



This drawing indicates that no terms of the sequence lie in the interval around 0.75; in particular, no tail lies in that interval. Finding an interval around a number  $\ell$  that contains no tail of the sequence is the same as finding an  $\varepsilon$  for which the 'approaching' property fails. So we can use drawings of sequences to figure out a good choice of  $\varepsilon$  to use in our proofs.

- Exercise 4.8.**
1. Consider the sequence  $A = (a_n = (-1)^n | n \in \mathbb{N})$ . Show that  $A$  does not 'approach' 1 by finding a specific positive real number  $\varepsilon$  such that no tail of  $A$  lies within a distance  $\varepsilon$  from 1. Similarly, show that  $A$  does not approach  $-1$ ,  $0$ ,  $2$ , or  $-2$ .
  2. Consider the sequence  $B = (b_n = 2^n | n \in \mathbb{N})$ . Show that  $B$  does not 'approach' any number  $\ell$ . This means, for each value of  $\ell$ , find a specific positive real number  $\varepsilon$  such that no tail of the sequence is

within  $\varepsilon$  from  $\ell$ ; that is, for each value of  $\ell$  find a specific number  $\varepsilon > 0$  such that, for infinitely many terms  $a_k$ ,  $\|a_k - \ell\| > \varepsilon$ .

- 3\*. Consider the sequence  $C = (c_n = \sin(n) | n \in \mathbb{N})$ . Show that  $C$  does not ‘approach’ any number  $\ell$ . You will need to look in the Appendix to find the precise definition of the trigonometric functions to do this thoroughly, but even without that you should be able to find the appropriate  $\varepsilon$  by drawing a picture.

We noticed above that the smaller the value of  $\varepsilon$  we choose, representing a smaller allowable error, then the more confident we were that we’ve chosen the correct number  $\ell$ . Our provisional definition of ‘approaches’ was satisfied by sequences that our intuition said were approaching a fixed number, and it incorporated this idea that smaller allowable errors increase certainty. Furthermore, sequences that our intuition said should not be converging did not satisfy the provisional definition. This means that the provisional definition is a good definition.

*Definition.* A sequence  $S = (a_1, a_2, a_3, \dots)$  **converges to a real number  $\ell$**  if and only if for any  $\varepsilon > 0$  there exists a natural number  $M_\varepsilon$  such that, for every natural number  $k \geq M_\varepsilon$ ,  $\|a_k - \ell\| < \varepsilon$ . A sequence  $T = (b_n | n \in \mathbb{N})$  **converges** if there exists a real number  $\ell$  such that the sequence  $T$  converges to  $\ell$ .

It is perhaps useful to think of convergence in terms of the sequence being able to meet any challenge. If any challenge  $\varepsilon$  is proposed, after some finite number of terms in the sequence are ignored, the remaining tail of the sequence lies within  $\varepsilon$  of the limit. Of course, when experimenting, if we choose a small value for  $\varepsilon$ , we should be more certain that we’ve actually chosen the right number  $\ell$ . However, as we saw in Exercise 4.4, we can’t stop after any particular choice of  $\varepsilon$ .

Understanding the definition of convergence is tricky because the definition involves infinitely many conditions, namely a condition for each  $\varepsilon$ . For example, if a sequence converges to 3, we know that after some point in the sequence, all the terms lie within a distance of 1 from 3, namely in the interval  $(2, 4)$ ; perhaps all the terms after the first hundred terms do so. But we also know that eventually all the terms will lie within a distance of 0.1 from 3, namely in the interval  $(2.9, 3.1)$ ; perhaps all the terms after the first million terms do so. We also know that eventually all the terms in the sequence lie within a distance of 0.001 from 3, namely in the interval  $(2.999, 3.001)$ ; perhaps all the terms after the first trillion terms do so. To converge, infinitely many such statements must be true.

To develop some intuition about convergent sequences, let's first look at the examples in the previous exercises and establish which ones converge and which ones do not.

- Exercise 4.9.*
1. Show that the sequence  $(\frac{1}{n}|n \in \mathbb{N})$  converges to 0.
  2. Show that the sequence  $(1 - e^{-n}|n \in \mathbb{N})$  converges to 1.
  3. Show that the sequence  $(\frac{(-1)^n}{n^2}|n \in \mathbb{N})$  converges to 0.
  4. Show that the sequence  $((-1)^n|n \in \mathbb{N})$  does not converge (to any number).
  5. Show that the sequence  $(2^n|n \in \mathbb{N})$  does not converge.
  - 6\*. Show that the sequence  $(\sin(n)|n \in \mathbb{N})$  does not converge.

There is one more important part of the intuitive notion of 'approaches': a sequence can only approach one number. This uniqueness was not part of the definition of a convergent sequence, but fortunately it does follow from that definition.

*Theorem 4.10.* If the sequence  $(a_n|n \in \mathbb{N})$  converges, then it converges to a unique number.

This last theorem tells us that a convergent sequence approaches exactly one number, which we will call the *limit* of the sequence.

*Definition.* If the sequence  $(a_n|n \in \mathbb{N})$  converges to  $\ell$ , then we say that  $\ell$  is the **limit** of the sequence. In this situation we write

$$(a_n|n \in \mathbb{N}) \rightarrow \ell.$$

All of this discussion of convergent sequences was motivated by the natural way in which we want to predict the position of a paint brush by looking at its position at nearby times.

*Exercise 4.11.* Consider the sequence

$$S = (100, \frac{100}{4}, \frac{100}{16}, \dots) = (\frac{100}{4^{n-1}}|n \in \mathbb{N}).$$

Check that  $S$  converges to 0, as officers Isaac and Gottfried claimed.



### 4.3 Finding Limits

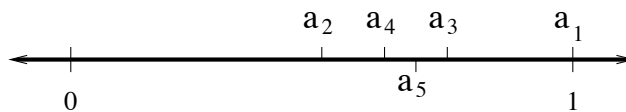
Officers Isaac and Gottfreid were very lucky that the limit of Zeno's position sequence was 0, because we needed to know to let  $\ell = 0$  to actually complete the proof that this sequence converged. It would be a little silly if we could only arrest people whose position at the moment of the crime was an integer, so we need to build better techniques for finding limits.

Thus far, every convergent sequence that we've considered has had an obvious limit. The sequence  $(\frac{100}{4^{n-1}} | n \in \mathbb{N})$  obviously converged to  $\ell = 0$ ; the sequence  $(1 - e^{-n} | n \in \mathbb{N})$  obviously converged to  $\ell = 1$ . But what about the sequence

$$S = (1, 1 - \frac{1}{2}, 1 - \frac{1}{2} + \frac{1}{4}, 1 - \frac{1}{2} + \frac{1}{4} - \frac{1}{8}, 1 - \frac{1}{2} + \frac{1}{4} - \frac{1}{8} + \frac{1}{16}, \dots)$$

$$= (\sum_{k=1}^n \frac{1}{(-2)^{k-1}} | n \in \mathbb{N}) = (1, 0.5, 0.75, 0.625, 0.6875, \dots)?$$

It may not be obvious looking at the numbers that this sequence converges, but we can think of the same sequence geometrically as follows. Imagine a person standing on the numberline at 0. He takes a step 1 unit to the right and then writes down his position. Then he takes a step 0.5 units to the left and writes down his position. Then he takes a step 0.25 units to the right and writes down his position. Repeating this process produces the sequence  $S$  above. It's intuitively obvious that the position of our person converges because he is alternating moving left-right and his steps keep getting smaller. If you think back to your last calculus class, you will remember this fact as the "alternating series test".



So we believe that this sequence converges, but what is the limit? Well, some of you may remember your calculus really well and may have noticed that  $S$  is the partial sums of a geometric series, so you have a formula that tells you the limit. For an example that you really don't have a limit for, consider the sequence

$$T = (1, 1 + \frac{1}{4}, 1 + \frac{1}{4} + \frac{1}{9}, 1 + \frac{1}{4} + \frac{1}{9} + \frac{1}{16}, 1 + \frac{1}{4} + \frac{1}{9} + \frac{1}{16} + \frac{1}{25}, \dots)$$

$$= \left( \sum_{k=1}^n \frac{1}{k^2} \right) = (1, 1.25, 1.36\bar{1}, 1.4236\bar{1}, 1.4636\bar{1}, \dots),$$

which converges by the “ $p$ -series test” to some number shrouded in mystery.

There are two directions that we can go from this impasse. First, we could try to find new ways to guess the limit of an arbitrary sequence. This is almost totally impossible, since every real number is the limit of a sequence, and very few of them have any special description in terms of the standard objects that we like (integers, fractions, roots,  $\pi$ , etc.). In fact, by using the decimal representation of a real number, we can realize it as the limit of a non-constant sequence. For example

$$(3, 3.1, 3.14, 3.141, 3.1415, 3.14159, \dots) \rightarrow \pi.$$

*Exercise 4.12.* Let  $R$  be a real number in the interval  $(0, 1)$ . Then  $R$  has an expression as a decimal number  $R = 0.R_1R_2R_3R_4\dots$  where the  $R_i$  represent the digits of the decimal expansion. Prove that the sequence

$$\tilde{R} = (0.R_1, 0.R_1R_2, 0.R_1R_2R_3, \dots)$$

converges to  $R$ . How would you generalize this so that  $R$  doesn’t have to be between 0 and 1?

*Corollary 4.13.* Every real number is the limit of a sequence whose terms are all rational numbers.

Corollary 4.13 shows that there are just too many sequences and limits to really hope to have a set of techniques that help us describe the limit of an arbitrary convergent sequence. So, instead, we could search for techniques to prove that a limit exists without knowing what that limit is. We will do this by finding properties of the sequences that don’t refer to the specific values of the terms in the sequences; instead the properties will describe the values of all terms at the same time or the relationships among the values.

Thus far, every sequence that we’ve worked with has been given by a formula. It’s great when this happens, but this is too much to ask in some situations. As a bonus, this new plan of attack will hopefully allow us to discuss the convergence of a sequence that isn’t given to us by a nice formula.

So let’s take a closer look at some of the convergent sequences we’ve seen

$$S_1 = \left( \frac{100}{4^{n-1}} \mid n \in \mathbb{N} \right) = (100, 25, 6.25, 1.57, 0.39, \dots)$$

$$S_2 = \left( \frac{1}{n} \mid n \in \mathbb{N} \right) = \left( 1, \frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \frac{1}{5}, \frac{1}{6}, \dots \right)$$

$$S_3 = (1 - e^{-n} | n \in \mathbb{N}) = (0.632, 0.865, 0.950, 0.982, \dots)$$

$$S_4 = \left(\frac{(-1)^n}{n^2} | n \in \mathbb{N}\right) = \left(-1, \frac{1}{4}, -\frac{1}{9}, \frac{1}{16}, -\frac{1}{25}, \dots\right)$$

and compare them to some sequences that don't converge.

$$S_5 = ((-1)^n | n \in \mathbb{N}) = (-1, 1, -1, 1, -1, 1, -1, 1, \dots)$$

$$S_6 = (2^n | n \in \mathbb{N}) = (2, 4, 8, 16, 32, 64, 128, \dots)$$

Then we will try to pick out some special properties that informed our intuition. Fortunately, we noticed one property long ago: each term in  $S_1$  is smaller than the previous term. Similarly, each term in  $S_2$  is also smaller than its predecessor, and each term in  $S_3$  is larger than the previous term. Let's give these two (related) properties names and precise definitions.

*Definition.* A sequence  $(a_n | n \in \mathbb{N})$  is called **increasing** if for any  $j < k$ ,  $a_j \leq a_k$ . The sequence is called **decreasing** if for any  $j < k$ ,  $a_j \geq a_k$ . A sequence is called **monotonic** if it is either increasing or decreasing. Note that a constant sequence,  $S = (c, c, c, \dots)$ , is both increasing and decreasing. If you are drawing this sequence on the real line, then it is monotonic if and only if it only moves in one direction (which includes the possibility of sometimes staying still as well).

Clearly being monotonic is not the same as convergent; the sequence  $S_4$  is not monotonic, but it does converge. Furthermore, the sequence  $S_6 = (2^n | n \in \mathbb{N})$  is increasing but does not converge. However, the notion of a sequence being monotonic does have the flavor that we were looking for; a sequence can be described as monotonic without ever giving a single specific term from the sequence. Instead, monotonicity describes the relationship among the terms of the sequence.

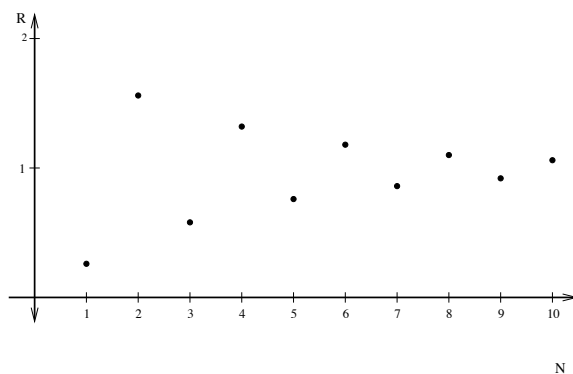
A second property of the sequence from Zeno's arrest jumps out as well: the terms in the sequence are all bigger than 0 (in fact, they are all between 0 and 100). This *bounded* property is a property that all four of the convergent sequences,  $S_1$  through  $S_4$ , share. Moreover,  $S_6$  does not have this property. It sounds useful, so let's give a formal definition.

*Definitions.* A sequence  $(a_n | n \in \mathbb{N})$  is **bounded from above** if there is a real number  $A$  such that  $a_k \leq A$  for any  $k \in \mathbb{N}$ . Similarly, the sequence is **bounded from below** if there is a real number  $B$  such that  $B \leq a_k$  for any  $k \in \mathbb{N}$ . A sequence is called **bounded** if it is bounded from below and bounded from above; equivalently, a sequence is bounded if there is a real number  $C$  such that for all  $k \in \mathbb{N}$ ,  $\|a_k\| \leq C$ .

*Exercise 4.14.* Take the six sequences above,  $S_1$  through  $S_6$ , and make a chart that includes one column for the sequence, one to say if it is monotonic or not, one to say if it is bounded or not, and one to say if it converges or not. Use this chart to make *several* conjectures about the relationships between the conditions of being monotonic, bounded, and convergent. (Not all of the conjectures need to involve all three ideas.)

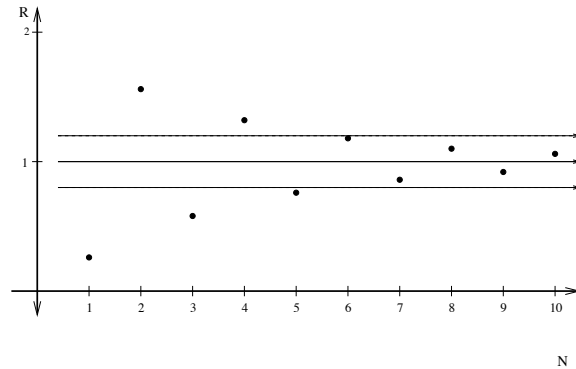
One useful way to make conjectures is to find a new way to represent the object you are studying and see if the new representation brings to light any new observations or sheds light on the relationships among the existing aspects you were considering. Thus far we have represented sequences as lists, as formulas, and as a bunch of marks on the real line. This last representation has the most obvious geometric uses, but it was messy. Perhaps we can find another graphical representation that doesn't have this problem. Fortunately, we all learned such a technique years ago: graphing. We're used to graphing functions like  $f(x) = 3\sqrt{x+5}$ , whose domain and codomain are subsets of  $\mathbb{R}$ . A sequence is a function from  $\mathbb{N}$  to  $\mathbb{R}$ ; for each natural number, the sequence gives is a real number.

Consider the sequence  $S = (a_n = 1 + (-\frac{3}{4})^n | n \in \mathbb{N})$ ; we could graph  $S$  as follows.



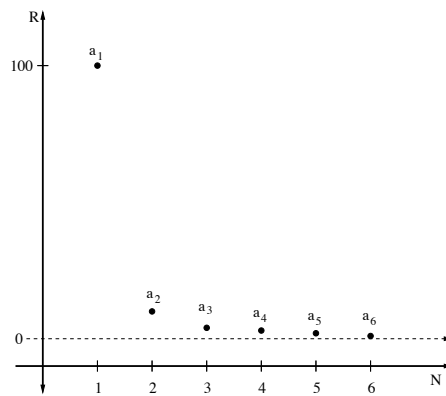
This sequence clearly will converge to  $\ell = 1$ , but what does that mean in terms of this graphical representation? Well, rather than being a point in one copy of  $\mathbb{R}$ ,  $\ell$  is a horizontal line in this new picture. And eventually being within  $\varepsilon$  from  $\ell$  means that to the right of some point, all points on the graph are inside an  $\varepsilon$ -tube of the line representing  $\ell$ . For example, this sequence is eventually within 0.2 from  $\ell = 1$ .

When drawing this picture, we are thinking of a sequence  $S$  as a function from  $\mathbb{N}$  to  $\mathbb{R}$ . Then this new representation of  $S$  is just the graph of this function. The idea that we can treat a sequence as a function will return in



the future sections.

Now let us use this new representation method to investigate our motivating examples. We start by drawing the representation of  $S_1$  and including the proposed limit.



Looking at the graph of  $S_1$ , the values approach the limit like the graph of a function approaches its horizontal asymptote. As we have noted, the sequence  $S_1$  from Zeno's graffiti arrest is decreasing and bounded. Since the sequence is decreasing, any limit must be smaller than all of the terms (namely a lower bound), but if the lower bound were too low, the graph would always stay far away from it. Our intuition used this information to guess the limit, which is the value of this asymptote. Put together, in relationship to  $S_1$ , we can describe the number 0 as the *greatest lower bound* for the sequence.

It turns out that the existence of a greatest lower bound for a subset of  $\mathbb{R}$  is actually a very subtle property, and we will simply state it as an axiom (a fact that we may use without justification). When looking at a subset of

$\mathbb{R}$  drawn as a line, it is easy to point at the greatest lower bound; however, it is not immediately clear whether you are pointing at a number or a hole between numbers. This axiom essentially claims that there are no “holes” in  $\mathbb{R}$ .

*Axiom* (Greatest Lower Bound Axiom). Let  $S \neq \emptyset$  be a subset of the real numbers with a real number  $L$  such that, for every  $s \in S$ ,  $L \leq s$ . Then there exists a *greatest lower bound* for  $S$ ,  $\inf(S)$ , called the **infimum** of  $S$ , with the following properties:

1. If  $s \in S$ , then  $\inf(S) \leq s$ .
2. If  $B$  is a real number such that  $B \leq s$  for all  $s \in S$ , then  $B \leq \inf(S)$ .

Note that the greatest lower bound axiom guarantees the existence of a similar, upper bound. This *least upper bound* is called the **supremum**, written  $\sup(S)$ . Furthermore, the properties in the axiom imply that the infimum and supremum of a set are unique.

*Exercise 4.15.* Write out a careful definition of the least upper bound axiom.

*Exercise 4.16.* Check that the infimum of a set is unique, if it exists. Carefully use only the properties guaranteed by the axiom, not your intuitive understanding of what the words should mean.

*Exercise 4.17.* For each of the following subsets of  $\mathbb{R}$ , argue whether or not the set has an infimum and supremum. Compute the infima and suprema that exist and justify your computations.

1.  $\mathbb{Q}$
2.  $(2, 5) \cup \{17\}$
3.  $\{\frac{1}{n} | n \in \mathbb{N}\}$

This axiom is useful because it allows us to produce a single number out of a complicated set. That sounds a lot like the process of finding a limit. Indeed, the axiom allows us to prove some new facts including our conjectures from Exercise 4.14. As we have seen, to prove that a sequence converges to  $\ell$  you need to be able to describe  $\ell$ . You have to know what  $\ell$  should be to compute  $\|a_k - \ell\|$  and compare it to  $\varepsilon$ . Neither of our abstract hypotheses in the conjectures (bounded/monotonic) produces the value of  $\ell$ , but the axiom will.

*Theorem 4.18.* Let  $S = (a_n | n \in \mathbb{N})$  be a sequence. If  $S$  is increasing and bounded from above, then  $S$  converges to  $\ell = \sup(\{a_n | n \in \mathbb{N}\})$ . Instead, if  $S$  is decreasing and bounded from below then  $S$  converges to  $\ell = \inf(\{a_n | n \in \mathbb{N}\})$ . In other words, bounded monotonic sequences converge.

This is a fabulous accomplishment, and we should be proud. But very few sequences are actually monotonic. Furthermore, there are still many sequences that speak to our intuition that have not influenced our theorems. For example, the following three sequences all obviously converge to 0.

$$A = (1, \frac{1}{2}, \frac{1}{3}, 7, 9, \frac{1}{4}, \frac{1}{5}, \frac{1}{6}, \frac{1}{7}, \frac{1}{8}, \frac{1}{9}, \dots)$$

$$B = (1, 0, \frac{1}{2}, 0, \frac{1}{3}, 0, \frac{1}{4}, 0, \frac{1}{5}, 0, \frac{1}{6}, 0, \dots)$$

$$C = (1, -1, \frac{1}{2}, -\frac{1}{2}, \frac{1}{3}, -\frac{1}{3}, \frac{1}{4}, -\frac{1}{4}, \dots)$$

None of these sequences is monotonic, and yet convergence seems obvious. This is because each of these sequences contains, buried inside it, a bounded monotonic sequence that we know converges, and our brains are able to filter the excess information. But this is not enough, it is also clear that all (but a finite number) of the other terms in the sequence are getting close to the same value as the buried sequence. Let's give a formal definition of this 'buried' sequence.

*Definition.* Let  $S = (s_n | n \in \mathbb{N})$  be a sequence. Then a **subsequence**,  $T$ , of  $S$  is a sequence obtained from  $S$  by omitting some of the terms of  $S$  while retaining the order. So  $T = (t_k = s_{n_k} | k \in \mathbb{N})$  subject to the condition that if  $i < j$ , then  $n_i < n_j$ .

This definition is really hard to parse; in particular, the subscript with its own subscript can be bewildering. So let's do an example. Consider the sequence

$$S = (s_n = 1 + 3n | n \in \mathbb{N}) = (4, 7, 10, 13, 16, 19, 22, 25, 28, 31, \dots),$$

which has the following two subsequences (among many others):

$$T = (t_k | k \in \mathbb{N}) = (4, 10, 16, 22, 28, 34, \dots) \text{ and,}$$

$$U = (u_k | k \in \mathbb{N}) = (7, 10, 25, 28, 31, 34, 37, \dots).$$

We could describe  $T$  as the sequence containing the odd terms (first, third, fifth...) from  $S$  (in the same order). So the first term of  $T$  is the first

term of  $S$ ; the second term of  $T$  is the third term of  $S$ , and similarly, the third term of  $T$  is the fifth term of  $S$ . In symbols,  $t_1 = s_1$ ,  $t_2 = s_3$ ,  $t_3 = s_5$ , and so forth. In particular,  $t_k = s_{2k-1}$ , so  $n_k = 2k - 1$ . It's easy to check that, if  $i < j$ , then  $n_i = 2i - 1 < 2j - 1 = n_j$ .

We could describe the subsequence  $U$  as the sequence formed from  $S$  by dropping the first and fourth through seventh terms. This is a much less regular pattern, so it will be harder to find a formula for  $n_k$ . But that doesn't mean we can't do it. Well,  $u_1 = s_2$ ,  $u_2 = s_3$ ,  $u_3 = s_8$ ,  $u_4 = s_9$ ,  $u_5 = s_{10}$ , and so on. So  $n_1 = 2$ ,  $n_2 = 3$ ,  $n_3 = 8$ ,  $n_4 = 9$ ,  $n_5 = 10$ , and so forth. We could collect this information as follows.

$$n_k = \begin{cases} k + 1 & \text{if } k \in \{1, 2\} \\ k + 5 & \text{if } k \geq 3 \end{cases}$$

Although it is important to understand this definition, most of the time it will be obvious from the description that the terms of a potential subsequence are in the same order as in the parent sequence.

*Exercise 4.19.* Let  $A$ ,  $B$ , and  $C$  be the sequences from above:

$$A = (1, \frac{1}{2}, \frac{1}{3}, 7, 9, \frac{1}{4}, \frac{1}{5}, \frac{1}{6}, \frac{1}{7}, \frac{1}{8}, \frac{1}{9}, \dots)$$

$$B = (1, 0, \frac{1}{2}, 0, \frac{1}{3}, 0, \frac{1}{4}, 0, \frac{1}{5}, 0, \frac{1}{6}, 0, \dots)$$

$$C = (1, -1, \frac{1}{2}, -\frac{1}{2}, \frac{1}{3}, -\frac{1}{3}, \frac{1}{4}, -\frac{1}{4}, \dots).$$

For each sequence, find a monotonic subsequence and describe  $n_k$  explicitly, checking that  $n_k$  increases with  $k$ .

Omitting terms from a sequence to produce a subsequence moves all subsequent terms towards the beginning of the sequence. This is a crucial property to use when working with subsequences, but it's hard to use in this form. So we state this more usefully in this lemma.

*Lemma 4.20.* Let  $S = (s_n | n \in \mathbb{N})$  be a sequence and  $T = (t_k = s_{n_k} | k \in \mathbb{N})$  a subsequence of  $S$ . Then, for every  $k \in \mathbb{N}$ ,  $k \leq n_k$ .

Some properties of parent sequences are not inherited by their subsequences.

*Exercise 4.21.* 1. Find a sequence that is not bounded with a bounded subsequence.

2. Find a sequence that is not monotonic with a monotonic subsequence.



3. Find a sequence that does not converge with a subsequence that does converge. Can you find a single sequence that works for all three parts of this problem? Can you find a single *formula* for a sequence that works?

Other properties of parent sequences are inherited by their subsequences.

*Theorem 4.22.* Let  $S$  be a bounded sequence and  $S'$  a subsequence of  $S$ ; then  $S'$  is bounded. Similarly, let  $T$  be a monotonic sequence and  $T'$  a subsequence of  $T$ ; then  $T'$  is monotonic.

*Theorem 4.23.* If the sequence  $S = (a_n | n \in \mathbb{N})$  converges to  $\ell$  and  $S'$  is a subsequence of  $S$ , then  $S'$  converges to  $\ell$ .

This last theorem gives us hope that we'll be able to find the limit of an arbitrary convergent sequence. If the limit exists, it is the same for a sequence and its subsequences. So we can use convergent subsequences to propose possible values for the limit. And as we all learned in highschool, knowing  $(\ell)$  is half the battle. The following is the key technical lemma in this battle strategy. It is hard, and you should try to draw several pictures using both visual representations of sequences to help outline your proof.

*Lemma 4.24.* Every sequence has a monotonic subsequence.

*Corollary 4.25.* Every bounded sequence has a convergent subsequence.

As mentioned above, having a convergent subsequence is not enough; the other terms in the sequence must get close to the values in the subsequence and the limit of that subsequence. Fortunately, we know the triangle inequality.

**Observation 2:** If two real numbers,  $x$  and  $y$ , are both 'very close' to  $\ell$ , then they are 'close' to each other. In symbols, if  $\|y - \ell\| < \varepsilon$  and  $\|x - \ell\| < \varepsilon$ , then  $\|y - x\| \leq \|y - \ell\| + \|\ell - x\| < 2\varepsilon$ . If a sequence  $S = (a_n | n \in \mathbb{N})$  converges to  $\ell$ , then after some point the terms in  $S$  are all 'very close' to  $\ell$ . If  $a_N$  and  $a_{N'}$  are two values that are within  $\varepsilon$  of  $\ell$ , then  $\|a_N - a_{N'}\| \leq \|a_N - \ell\| + \|\ell - a_{N'}\| < 2\varepsilon$ .

In other words, using the triangle inequality, we see that the condition of all terms in a sequence being close to a value  $\ell$  is related to the condition that all terms are close to each other. The notion that eventually all of the terms are very close to each other is another property that doesn't refer to a specific term's value. So we give it a precise definition.

*Definition.* A sequence  $(a_n | n \in \mathbb{N})$  is called a **Cauchy sequence** if for every  $\varepsilon > 0$  there exists an  $N_\varepsilon$  such that for all  $i, j \geq N_\varepsilon$ ,  $\|a_i - a_j\| < \varepsilon$ .

Note that this definition does not say that the distance between *consecutive* terms is small; it says something much stronger. It says that, after some point, the distance between a term and all subsequent terms must be small.

*Exercise 4.26.* Find an example of a sequence such that the distance between consecutive terms decreases to 0 but the sequence does not converge. (Hint: Consider the harmonic series that you studied in calculus.)

*Exercise 4.27.* Check directly (not as a corollary) that the sequence

$$(3, 2.1, 2.01, 2.001, \dots) = (2 + (\frac{1}{10})^{n-1} | n \in \mathbb{N})$$

is a Cauchy sequence.

*Theorem 4.28.* Suppose the sequence  $S = (a_n | n \in \mathbb{N})$  converges, then  $S$  is bounded. Similarly, suppose the sequence  $T = (b_n | n \in \mathbb{N})$  is Cauchy, then  $T$  is bounded.

*Exercise 4.29.* In Theorem 4.28 you proved that convergent sequences are bounded and that Cauchy sequences are bounded. Compare your proofs of these two facts.

As Observation 2 points out, convergent sequences are Cauchy sequences, but you should prove this fact carefully. The issue to address is the unwanted 2 that appeared in the observation.

*Theorem 4.30.* Let  $S$  be a convergent sequence; then  $S$  is a Cauchy sequence.

Thus far, it seems that we've found yet another property satisfied by convergent sequences. But if you think back carefully to the proofs of our theorems, most of them can easily be adopted from the case of convergent sequences to the case of Cauchy sequences using some simple triangle inequality tricks as in the proof of Theorem 4.30. For example, prove the following theorem under the two different hypotheses. (Of course you could prove the second sentence as a corollary of the first using Theorem 4.30, but that would defeat the purpose.)

Depending on how your sequence is presented to you, it may be easy or difficult to check that it is Cauchy; however, this check is certainly at worst as difficult as checking that your sequence converges. Moreover, the property of being Cauchy does not require knowledge of an elusive  $\ell$  and is intrinsic to the sequence, meaning that it only uses the terms in the sequence. It turns out that the Cauchy and convergent properties are equivalent. So the property of being a Cauchy sequence is the answer to our problems with the definition of a convergent sequence.

*Theorem 4.31.* Let  $S$  be a sequence. Then  $S$  is convergent if and only if  $S$  is Cauchy.

This last theorem tells us that we can take the notion of a Cauchy sequence as the definition of convergence, and it is preferable because the Cauchy condition is intrinsic to the sequence.

Thus far in this chapter, we have been dealing with the real numbers, but the notion of a Cauchy sequence makes sense much more broadly. Furthermore, we can use the rational numbers,  $\mathbb{Q}$ , and the notion of a Cauchy sequence to *define* the real numbers, a feat that we are totally ignoring in this course. In fact, this is essentially what the decimal representation of a real number is.

## 4.4 Continuity

It seems like the case has been nailed pretty firmly shut around Zeno. But the rabbit hole goes deeper. The reason we can't lock Zeno up right now is that officers Isaac and Gottfried *chose* to measure Zeno's position using time intervals of length  $(\frac{1}{2})^n$  minutes. A careful judge would find this evidence circumstantial. It is conceivable that if the officers had used time intervals of length  $(\frac{1}{10})^n$  minutes, then Zeno would not have appeared to be approaching the wall. Even one such piece of evidence would have Zeno appealing his conviction immediately.

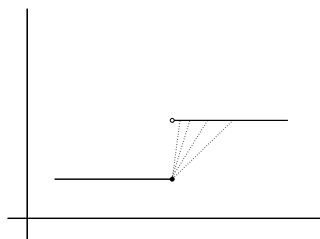
The thing that would convince such a discerning judge is knowing that *any* choice of time intervals would produce the same value for Zeno's position at 1am. So to complete our case we need to show that for any choice of time intervals, the sequence of positions converges and each of these sequences converges to the same value.

This sounds like a daunting amount of information that we must know about Zeno's position: to build a really strong argument against Zeno, we need to know his position at *any* point in time before 1am (at least near the time in question), which is the same thing as knowing his position as a *function* of time. This function takes input in the form of an amount of time (a real number) and produces a position (a real number distance from the wall), thus Zeno's position is a function from the real numbers to the real numbers,  $p : \mathbb{R} \rightarrow \mathbb{R}$ . So now we are going to turn our attention to **functions on the reals**, functions whose domains and codomains are (subsets of)  $\mathbb{R}$ .

Most of the time, we will need a formula for a function on the reals to actually complete these computations, like Zeno's position function. We are

all familiar with many functions from the real numbers to the real numbers such as  $f(x) = x^3$ ,  $g(x) = \sin(x)$ , and  $h(x) = e^x$ . However, functions on the reals need not have a neat expression. For example, the function that sends each rational number to 0 and each irrational number to 1 is a perfectly good function on the reals.

In the next few sections we are hoping to show when the choices made while computing the position of Zeno's brush at 1am do not effect the value obtained for the answer. This is a property of the function that describes Zeno's position, and not all functions have this property. For example, if his brush were able to teleport (meaning move some distance in absolutely no time), then there is no way that these computations could produce anything meaningful. The following graph represents the position of a brush that is not moving, suddenly teleports, and then remains still at the new position. You can see that sequence of positions is constant, but the position at moment of teleportation is not the limit of the times leading up to that moment.



So we must assert that Zeno's brush cannot teleport and further investigate what properties we are assuming Zeno's position function satisfies. But first, we have been using a procedure for producing new sequences, given a sequence and a function. Let's make sure that it is clear.

*Definition.* Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a function and let  $S = (a_n | n \in \mathbb{N})$  be a sequence. Then define the sequence  $f(S) = (f(a_n) | n \in \mathbb{N})$ , which we will sometimes call the **image sequence** or the image of  $S$  under  $f$ .

*Exercise 4.32.* Let  $f(x) = 2x + 1$ .

1. Consider the sequence  $T = (\frac{1}{n} | n \in \mathbb{N})$ . Compute the sequence  $f(T)$  and show that it converges to  $f(0) = 1$ .
2. Similarly, consider the sequence  $U = ((\frac{-1}{10})^{n-1} | n \in \mathbb{N})$ . Compute the sequence  $f(U)$  and show that it converges to  $f(0) = 1$ .

*Exercise 4.33.* Let  $F : \mathbb{R} \rightarrow \mathbb{R}$  be the function defined by

$$F(x) = \begin{cases} x & \text{if } x \in \mathbb{Q} \\ 1 & \text{if } x \in \mathbb{R} \setminus \mathbb{Q} \end{cases}.$$

1. Let  $S = (\frac{1}{n} | n \in \mathbb{N})$ . Compute the sequence  $F(S)$  and find its limit.
2. Let  $T = (\frac{\sqrt{2}}{n} | n \in \mathbb{N})$ . Compute the sequence  $F(T)$  and find its limit.

The previous exercise has shown us that there is something funny with this special function  $F$  that produced two different values for the limit of the image sequence. In particular, both  $S$  and  $T$  converge (in the domain of  $F$ ) to 0, but  $F(S)$  and  $F(T)$  converge to different values in the codomain. Even worse, there are sequences whose image doesn't even converge.

*Exercise 4.34.* Let  $F(x)$  be the function defined in Exercise 4.33 that sends each rational number to itself and each irrational number to 1. Find a convergent sequence  $S$  such that the sequence  $F(S)$  is not convergent.

Fortunately, this does not happen with the function that describes the position of Zeno's car or most of the other functions that you're familiar with from calculus.

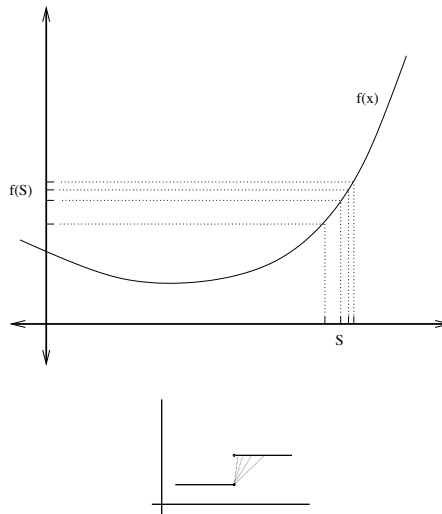
*Exercise 4.35.* 1. Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be the function  $f(x) = 2x + 1$ . Let  $S = (a_n | n \in \mathbb{N})$  be a convergent sequence, and show that the sequence  $f(S) = (2(a_n) + 1 | n \in \mathbb{N})$  is a convergent sequence. (Hint: If  $S$  converges to  $\ell$ , then the limit of  $f(S)$  should be  $f(\ell)$ .)

- 2\*. Let  $g : \mathbb{R} \rightarrow \mathbb{R}$  be the function  $g(x) = \cos(x)$ . Let  $S = (a_n | n \in \mathbb{N})$  be a convergent sequence, and show that the sequence  $g(S) = (\cos(a_n) | n \in \mathbb{N})$  is a convergent sequence.

To make sense of the instantaneous velocity of an object whose position is described by the function  $f : \mathbb{R} \rightarrow \mathbb{R}$ , we need to rule out cases like the teleportation example. Thus far, the functions that seem to produce useful computations send convergent sequences to convergent sequences with the appropriate limit. So we give a name to functions with this property.

*Definition (Limit Continuity).* A function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is **continuous** if and only if for every sequence  $S = (a_n | n \in \mathbb{N})$  that converges, with limit  $\ell_S$ , the sequence  $f(S) = (f(a_n) | n \in \mathbb{N})$  converges to  $f(\ell_S)$ .

As before, this definition is somewhat difficult to understand. Here is an example of the graph of a function  $f$  along with a sequence  $S$  converging in its domain where  $f(S)$  converges in the codomain (to the appropriate limit).



And here is an example of the graph of a function  $g$  along with a sequence  $T$  converging to  $\ell$  in its domain such that  $f(T)$  does not converge to  $g(\ell)$ .

Notice that the function,  $g$ , which is not continuous, has a gap in its graph. This is not a rigorous notion, so you should not use it in your proofs. But you can use the idea that continuous functions have graphs that can be drawn without lifting the pencil to inform your intuition.

Many classes of functions are continuous.

*Theorem 4.36.* For any real numbers  $a$  and  $b$ , the function  $f : \mathbb{R} \rightarrow \mathbb{R}$  defined by  $f(x) = ax + b$  is continuous.

*Theorem 4.37.* The function  $f : \mathbb{R} \rightarrow \mathbb{R}$  defined by  $f(x) = \|x\|$  is continuous.

*Theorem\* 4.38.* The trigonometric functions  $\sin(x)$  and  $\cos(x)$  are continuous.

*Theorem\* 4.39.* The exponential function  $e^x$  is continuous.

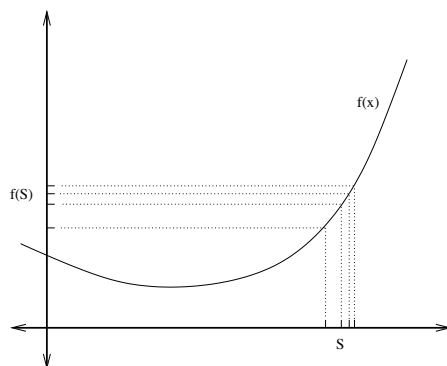
Checking that a function is (limit) continuous involves checking that a large collection of sequences converge. For some functions, all sequences can be dealt with simultaneously, as you probably did with the linear functions in the proof of Theorem 4.36. For other functions, the proof has a different flavor depending on the value of  $\ell$ , the limit point in the domain, as you probably did with the absolute value function in the proof of Theorem 4.37. The fact that we can check continuity by checking it at each point in the domain means that continuity is what's called a *local* condition.

You may also have noticed that the sequences actually play a relatively small role in the proofs of continuity. You must prove that *every* sequence

behaves in the same way, so you may make no interesting assumptions about your sequences. The definition is saying something about every tail of every sequence converging to a fixed point  $x$  in the domain of  $f$ . This set sounds very complicated, but it is in fact quite simple.

Observation 3: If a sequence converges to  $x$ , then for any distance  $\delta$  there is a tail within  $\delta$  for  $x$ , namely, there is a tail contained in the interval  $(x - \delta, x + \delta)$ . Conversely, every number in the interval  $(x - \delta, x + \delta)$  is in a tail of a sequence converging to  $x$ .

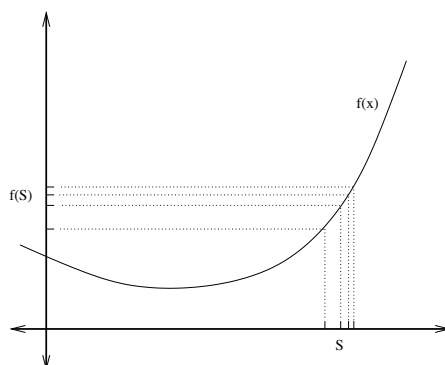
Observation 3 suggests that rather than formulating continuity in terms of all sequences, it can be formulated in terms of an interval  $(x - \delta, x + \delta)$ . But exactly what should be the definition? Remember that the definition of continuity requires the image sequences to converge to the correct point in the codomain. This means that, for each  $\varepsilon$ , the image of the tails must be within  $\varepsilon$  of the proposed limit. In our graphical representation, as above, being within  $\varepsilon$  of  $f(x)$  is most naturally represented by the graph of the function being within the  $\varepsilon$ -tube of  $f(x)$ .



But *what* must be inside the  $\varepsilon$ -tube? For every  $\varepsilon$ , there must be a tail, and it is the interval  $(x - \delta, x + \delta)$  that represents the tails. So, if  $f$  is continuous, for every  $\varepsilon$  in the codomain around  $f(x)$ , there is  $\delta$  such that  $(x - \delta, x + \delta)$  maps into  $(f(x) - \varepsilon, f(x) + \varepsilon)$ .

*Definition.* A function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is **continuous at a point**  $x_0 \in \mathbb{R}$  if and only if for every  $\varepsilon > 0$ , there exists a  $\delta > 0$  such that, for every  $y \in \mathbb{R}$  with  $\|y - x_0\| < \delta$ ,  $\|f(y) - f(x_0)\| < \varepsilon$ .

So being continuous at a point means that for every challenge,  $\varepsilon$ , there is a response to the challenge,  $\delta$ , such that points closer than  $\delta$  to  $x_0$  are



taken to points less than  $\varepsilon$  from  $f(x_0)$ . This definition tells us the meaning of being continuous at a point  $x_0$ . To be continuous overall, this property must occur at each point in the domain.

*Definition* (Metric Continuity). A function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is **continuous** if it is continuous at every point. Note that for each point,  $x$ , there is a  $\delta_x$  satisfying the inequalities above, but these  $\delta_x$  need not have anything to do with each other.

Since we have two definitions of continuity, we had better confirm that the two definitions are equivalent. It should be clear from this development why metric continuity implies limit continuity, but the converse is not obvious.

*Theorem 4.40.* The two definitions of continuity are equivalent to one another. In other words, a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is continuous in the limit sense if and only if it is continuous in the metric sense.

We proved, in the previous section, that every Cauchy sequence converges in  $\mathbb{R}$ , which is a formalization of the idea that  $\mathbb{R}$  has “no holes”. Requiring a function to be continuous is like asking it not to create any holes, that its graph can be drawn without “picking up the pencil”. This is yet another instance of one of our major themes in this book: study objects with some structure and the functions between them that preserve that structure.

The structure we are considering here is the notion of convergent sequences, which came from a sense of distance on  $\mathbb{R}$ . That is, if a sequence  $S$  is converging to a number  $\ell$ , then the terms of  $S$  are getting close to one another and to  $\ell$ . The definition of continuity tells us that the image sequence,  $f(S)$ , does the same thing, namely approaches  $f(\ell)$ .

If we instead considered the important structure as distances, then we



would be interested in functions that send points that are close to points that are close. In trying to define this property, we would have run into the old problem from the convergence section: how close is close enough? We could apply the same process of observations and refinements that we followed in the section on convergence to produce a good definition for this property. As before, we would attack the problem by claiming that we can do better than any given error,  $\varepsilon$ . More specifically, for a continuous function  $f : \mathbb{R} \rightarrow \mathbb{R}$  and a point  $x_0$ , given a distance  $\varepsilon$ , we find an interval  $(x_0 - \delta, x_0 + \delta)$  that maps into the interval  $(f(x_0) - \varepsilon, f(x_0) + \varepsilon)$ . This is exactly the  $\varepsilon$ - $\delta$  definition of continuity that we found by other means. Both definitions are standard, but the  $\varepsilon$ - $\delta$  definition is used more commonly when thinking about the theorems that we're proving.

For now, let's return to the position of Zeno's Mustang (mileage marker) as a function of time ( $t$  minutes after 3:00pm), and call this function  $p(t)$ . Remember that he drove so that at time  $t$  the nose of his car was at mileage marker  $t^2$ , so  $p(t) = t^2$ . Our next goal is to show that  $p$  is continuous. We could try to do this directly, or we could try to prove that a broad class of functions is continuous, encompassing the example of Zeno's case. Since it's simply fortuitous that Zeno's speed was such a simple function, we will take this latter approach.

Combining continuous functions through addition, multiplication, or composition yields continuous functions.

*Theorem 4.41.* Let  $f(x)$  and  $g(x)$  be continuous functions from  $\mathbb{R}$  to  $\mathbb{R}$ . Then  $(f + g)(x)$ , defined as  $(f + g)(x) = f(x) + g(x)$ , is continuous.

*Theorem 4.42.* Let  $f(x)$  and  $g(x)$  be continuous functions from  $\mathbb{R}$  to  $\mathbb{R}$ . Then  $(fg)(x)$ , defined as  $(fg)(x) = f(x)g(x)$ , is continuous.

*Theorem 4.43.* Let  $f(x)$  and  $g(x)$  be continuous functions from  $\mathbb{R}$  to  $\mathbb{R}$ . Then  $(g \circ f)(x)$ , defined as  $(g \circ f)(x) = g(f(x))$ , is continuous.

*Corollary 4.44.* Any polynomial  $p(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x^1 + a_0$  is continuous.

*Lemma 4.45.* Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be continuous at  $x_0$ . If  $f(x_0) > 0$ , then there is a  $\delta > 0$  such that, for every  $y \in \mathbb{R}$  with  $\|y - x_0\| < \delta$ ,  $f(y) > 0$ . Moreover, there is a  $\delta' > 0$  such that, for every  $y \in \mathbb{R}$  with  $\|y - x_0\| < \delta'$ ,  $f(y) > \frac{f(x_0)}{2}$ .

*Theorem 4.46.* Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a continuous function that is never 0. Then the function  $h : \mathbb{R} \rightarrow \mathbb{R}$  defined as  $h(x) = \frac{1}{f(x)}$  is also continuous.

These theorems allow us to show that vast numbers of functions are continuous, such as  $\sin(e^x)\tan(x^2 + 3x + 4)/(x^2 + 1)$ .

*Exercise 4.47.* Show that the position function for Zeno's car,  $p(t) = t^2$ , is continuous as a corollary and directly (without using any of these theorems). Could you extend your direct proof to work for an arbitrary polynomial?

We embarked on this investigation in part to describe a property that ruled out the possibility that Zeno's car was teleporting. We've said that continuous functions can be thought of, conceptually, as functions whose graphs we can draw without lifting the pencil. The theorem that really captures this sense is the Intermediate Value Theorem, which states that if a function takes on two values, then it must also take on every value in between. Notice that the following theorem is false for functions from  $\mathbb{Q}$  to  $\mathbb{Q}$ , so we must again use one of the axioms of the Reals somewhere in the proof.

*Theorem 4.48* (Intermediate Value Theorem). Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a continuous function and let  $a$  and  $b$  be two real numbers such that  $f(a) < f(b)$ . Then for any real number  $r$  such that  $f(a) < r < f(b)$ , there is a real number  $c \in [a, b]$  such that  $f(c) = r$ .

## 4.5 Speeding and Zeno's Paradox<sup>TM</sup>

Speeding tickets are the bane of existence for people who speed (which is almost everyone in Austin). So, when Zeno conceived of his patented Paradox<sup>TM</sup>, he was pretty sure he was about to retire to the lap of luxury. The great advantage of the Paradox<sup>TM</sup> product over those other (radar detecting) devices was that it did not involve slowing down! It all would have worked perfectly except that the cops who pulled him over during the test run were officers Isaac and Gottfried, whose mathematical insights had wrested order from the jaws of vehicular anarchy. But we have gotten ahead of ourselves; the story begins with Zeno putting his new Mustang to its paces.

One spring afternoon our "hero", Zeno, jumped into his Mustang convertible and galloped down the springtime highway. The "30 miles per hour" speed limit signs were a mere blur as he raced by. He kept his speed so that at  $m$  minutes after 3pm he was exactly at the mileage marker  $m^2$ . Soon the serenity of the sunny drive exploded as sirens blared, lights flashed, and the strong arms of the law pulled Zeno over for speeding. Zeno had talked his way out of tons of tickets in his life, and he felt his Paradox<sup>TM</sup> was easily up to the current challenge. So Zeno had no fear that the approaching officers would overcome his evidence of innocence. But his confidence might have been a bit shaken if he had noticed that the two officers who approached

his window really knew their math. The officers walked up to Zeno's rolled down window and asked:

Officer Gottfried: Do you know why I pulled you over, sir?

Zeno: No officer, I don't.

Officer Gottfried: Well, the speed limit is 30 miles per hour, and you were doing 120; that's two miles per minute!

Zeno: Really? When?

Officer Gottfried: At precisely 3:01pm.

Zeno: You must be mistaken. At 3:01pm, precisely, I was not moving at all, and I can prove it.

Officer Gottfried: How can you prove it?

Zeno: My Zeno's Paradox<sup>TM</sup> recorded the whole story. You will see that at precisely 3:01 I was only in one place. Here is an instant photograph supplied by the Paradox<sup>TM</sup> that shows explicitly where I was at precisely 3:01. The Paradox<sup>TM</sup> was cleverly located exactly across the street from the 1 mile mileage marker sign. And you see that at the exact moment, the nose of my Mustang is precisely lined up with the 1 mile marker. You see that the picture is time-stamped 3:01 exactly. Can I go now?

Officer Isaac: Hold your horses, Bud. You aren't the only cowboy with a camera. Here is a photo of you at 3:02 exactly with that Mustang nostril lined up at the 4 mile marker. Now if I know my math, and you'd better believe I do, that means you went 3 miles in 1 minute, which is why you are going down my friend.

Zeno: Put those cuffs away. You haven't made your case. The question is not where I was at 3:02, the question is how fast I was going at 3:01 and my snapshot shows I was in one place and I rest my case.

Officer Isaac: You will rest your case alright, and you'll rest it in the slammer, because we've got more evidence. Here's another picture—your Mustang's snozzola at the 1.21 mileage marker at precisely 3:01.1. So you were at the 1 mileage marker at 3:01 and the 1.21 mileage marker at 3:01.1. So you went 0.21 miles in .1 minutes. That works out to 2.1 miles per minute during that half a minute.

Zeno: I'm getting bored. What does my location at 3:01.1 have to do with the question at hand? We are supposed to be talking about 3:01, and at 3:01 I was in precisely one place.

Officer Gottfried: Unfortunately for you, we had an infinite number of cameras taking an infinite number of pictures. And, altogether, they tell a convincing story about speeding: At 3:01.01 you were at mileage marker 1.0201. That means that you went  $1.0201 - 1 = .0201$  miles in 0.01 minutes;

that is an average speed of 2.01 miles per minute. At 3:01.001, you were at mileage marker 1.002001. So you went  $1.002001 - 1 = 0.002001$  miles in 0.001 minutes. That is an average speed of 2.001 miles per minute. We noticed that your location at each time 3:01 plus  $\Delta$  minutes was exactly at mileage marker  $(1 + \Delta)^2$ . So for every interval of time  $\Delta$  after 3:01, your average speed during the interval of time from 3:01 until 3:01 +  $\Delta$  was  $\frac{(1+\Delta)^2-1}{\Delta}$  (which equals  $2 + \Delta$ ). You are right that no *one* piece of evidence is conclusive, but the totality of this infinite amount of evidence with arbitrarily small lengths of time tells the story. Your instantaneous velocity at time 3:01 was 2 miles per minute because your average velocities during tiny lengths of time around 3:01 *converge* to 2 miles per minute. Zeno your speeding days are done.

Zeno: Converge? What does “converge” mean? I admit it looks bad for me and my Paradox<sup>TM</sup>, but I’m not going to give up meekly until you explain precisely what you mean by *converge*.

Officer Gottfried: Okay, let’s just think it over for 150 years or so and then you will be convinced.

Officer Isaac: You should have plenty of time to ponder during your night in the slammer.

Zeno: Can we speed this up, I’ve got places to be?

---

The moral of this story is that speed is not a directly measurable quantity. We can measure the mass of Zeno’s car directly because all matter exerts a force on other matter, proportional to its mass. Ignoring relativity, we can even measure position with a ruler and time with a stop watch or clock. But to measure speed, we must measure *other* quantities (time and position) at least twice and compute an *average* speed. As the story above indicates, the closer together our measurements are spaced, the greater accuracy we have about what’s going on. But no car ever drives the same speed for any length of time; even using the cruise control, the car is changing speeds due to hills and other tiny factors. The solution seems to be that we should take *many* measurements of average speed. This story shows that instantaneous velocity requires us to make infinitely many average speed computations using elapsed times that get arbitrarily close to zero elapsed time, but zero elapsed time makes no sense for measuring motion. Zeno’s *instantaneous velocity* is the number to which the totality of an infinite number of computations of average speeds over progressively shorter intervals of time converges.

Specifically, if Zeno's position at every time  $t$  is given as a function  $p(t)$ , then the instantaneous velocity at any specific time  $t_0$  is the number to which the values  $\frac{p(t_0+\Delta)-p(t_0)}{\Delta}$  converge as we select values of  $\Delta$  that get close to 0. The process of finding the instantaneous velocity as a single number that summarizes all the approximations is called "taking the limit" or "converging to a limit". In our example, we got a sequence of average velocities computed over progressively shorter intervals of time. By looking at intervals of time of length 1 minute, then 0.1 minute, then 0.01 minute, then 0.001 minute, then 0.0001 minute (each interval starting at 3:01pm), we computed the average velocities to get a sequence of average velocities

$$(3, 2.1, 2.01, 2.001, 2.0001, \dots).$$

We plausibly concluded that this sequence of numbers converges to 2, therefore concluding that Zeno's instantaneous velocity at 3:01pm was 2 miles per minute. But now we are left with pinning down what *convergence* really means.

When putting together the case against Zeno, we considered the position of his car as a function of time. We then took a sequence of times (in the domain of this function) and put them into this function to produce a sequence of positions (in the codomain of the function). We then used this sequence of positions to produce a sequence of average velocities. Finally, we computed the limit of this last sequence and called it his instantaneous velocity.

$$(s_n)_{n \in \mathbb{N}} \rightsquigarrow \left( \frac{p(1+s_n) - p(1)}{(1+s_n) - 1} \right)_{n \in \mathbb{N}}$$

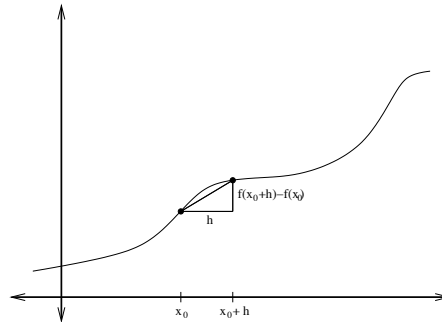
When computing the terms in this last sequence, we are repeatedly computing the fraction  $\frac{\text{distance travelled}}{\text{time elapsed}}$ . We will see this complicated fraction involved in computing average velocities many times in the future, so we give it a name.

*Definition.* Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a function. Let  $x_0$  be an number in the domain of  $f$  and define a new function  $\Delta(f, x_0) : \mathbb{R} \setminus \{0\} \rightarrow \mathbb{R}$  by

$$\Delta(f, x_0)(h) = \frac{f(x_0 + h) - f(x_0)}{(x_0 + h) - x_0} = \frac{f(x_0 + h) - f(x_0)}{h}$$

called the **difference quotient** of  $f$  at  $x_0$ . Note that this difference quotient is a function of  $h$ , corresponding to the elapsed time.

Graphically, the difference quotient of  $f$  at  $x_0$  evaluated at  $h$  is just the slope of a certain line.



## 4.6 Derivatives

Let's now return to our hapless "hero", Zeno, who is not doing well in his attempt to avoid his just punishment. Recall that the two officers Isaac (last name Newton) and Gottfried (last name Leibniz) had presented strong reasoning that Zeno's instantaneous velocity could be computed by taking the limit of  $\frac{p(t_0 + \Delta) - p(t_0)}{\Delta}$  as  $\Delta$  goes to 0. Now that we have a better idea of what a limit means, we realize that for some functions  $p(t)$ , the limit that must be taken to determine the instantaneous velocity may or may not exist. So let's first give a name to this process of computing the instantaneous velocity.

*Definition.* Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a function. If, for any sequence  $S$  converging to  $x_0$ ,  $f(S)$  converges, to the same point  $y$ , then we say that "the limit as  $x$  goes to  $x_0$  of  $f(x)$  is  $y$ " and write

$$\lim_{x \rightarrow x_0} f(x) = y.$$

If there is any sequence  $T$  converging to  $x_0$  such that  $f(T)$  does not converge or if there are two sequences  $S_1$  and  $S_2$  converging to  $x_0$  such that  $f(S_1)$  and  $f(S_2)$  converge to different values, then we say that  $\lim_{x \rightarrow x_0} f(x)$  does not exist.

*Definition.* Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a function. For any real number  $x_0$ , the derivative of  $f$  at  $x_0$ , denoted  $f'(x_0)$ , is  $\lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}$ , if that limit exists. When the limit does exist, we say that  $f(x)$  is **differentiable at  $x_0$** . If  $f(x)$  is differentiable at each point  $x$  in its domain, then  $f$  is **differentiable**.

Differentiable functions can be thought of conceptually as functions whose graphs at each point look straight when looked at under a high powered microscope. Can you see why the definition of the limit and the definition of derivative tell us that differentiable functions are ones that look straight

when magnified? The fact that differentiable functions locally look like a straight line means in particular that differentiable functions are continuous.

*Theorem 4.49.* A differentiable function is continuous.

We now proceed essentially to duplicate all our work on continuous functions, but this time considering differentiability instead of continuity. Many classes of functions are differentiable. Proving that functions are differentiable is generally more difficult than proving that they are continuous, because we need to prove that the more complicated difference quotient has a limit. The next several theorems allow us to prove that polynomials are differentiable.

When proving that a particular function is differentiable, we must always return to the definition of differentiability, namely, we must prove that  $\lim_{h \rightarrow 0} \frac{f(x_0+h)-f(x_0)}{h}$  exists. Establishing what the derivative of a function is means to give another function that equals that limit.

*Theorem 4.50 (Power Rule).* For every natural number  $n$ ,  $f(x) = x^n$  is differentiable and  $f'(x) = nx^{n-1}$ .

*Exercise 4.51.* Use the Power Rule to compute Zeno's instantaneous velocity at 3:01pm.

*Theorem 4.52.* If  $f(x)$  is differentiable and  $a$  is a real number, then the function  $g(x) = af(x)$  is also differentiable and  $g'(x) = af'(x)$ .

*Theorem 4.53.* Let  $f(x)$  and  $g(x)$  be differentiable functions from  $\mathbb{R}$  to  $\mathbb{R}$ . Then  $(f+g)(x)$ , defined as  $(f+g)(x) = f(x) + g(x)$ , is differentiable and  $(f+g)'(x) = f'(x) + g'(x)$ .

Before we get too lackadaisical about how these derivatives are going to proceed, let's point out that products do not work as expected.

*Exercise 4.54.* Find two differentiable functions  $f(x)$  and  $g(x)$  for which the derivative of their product is not the product of their derivatives.

We will see in a moment what the derivative of a product of two differentiable functions is, but first let's complete the polynomial story.

*Corollary 4.55.* Every polynomial function

$$f(x) = a_n x^n + a_{n-1} x^{n-1} + a_{n-2} x^{n-2} + \dots + a_2 x^2 + a_1 x + a_0$$

is differentiable and

$$f'(x) = na_n x^{n-1} + (n-1)a_{n-1} x^{n-2} + (n-2)a_{n-2} x^{n-3} + \dots + 2a_2 x + a_1.$$

In Exercise 4.54 above you saw that the derivative of a product is not quite as simple as one might think. So let's analyze what the derivative of the product is and why it is so.

Before we proceed with the derivative of the product of two functions, let's introduce an alternative notation for the derivative. If  $f(x)$  is a function then its derivative  $f'(x)$  can be denoted  $\frac{d}{dx}(f(x))$ . Notice that this notation reminds us of the definition of the derivative. This notation was carefully designed to do so by one of the inventors of calculus Gottfried Leibniz. Leibniz thought carefully about the notation so that operations of calculus could be done somewhat mechanically. One of the virtues of calculus is that much calculus work can be done by rote, and Leibniz's carefully crafted notation makes such routine work convenient.

Let's now return to analyzing the derivative of the product of two functions. We will begin by considering the product of two specific, simple functions.

*Exercise 4.56.* Let  $f(x) = 2x$  and  $g(x) = 3x$ . What is the derivative of the product, that is,  $\frac{d}{dx}(f(x)g(x))$ , at the point  $x = 4$ ? Of course, you could simply multiply  $f(x)$  times  $g(x)$  to get the product  $6x^2$  and then take its derivative. That is fine, particularly to check your thinking; however, for this exercise please think about the definition of the derivative. When you put in the  $h$  into the numerator of the definition, what happens? Try to follow the reasoning through, thinking about the derivatives of each of  $f(x)$  and  $g(x)$  as you do. The goal of this exercise is for you to see the relationships among the derivatives of each of the functions, the values of each of the functions, and the derivative of the product.

If you were successful with the previous exercise or if you remember the Product Rule from a calculus course, the following theorem will not be a surprise.

*Theorem 4.57 (Product Rule).* Let  $f(x)$  and  $g(x)$  be differentiable functions. Then their product is differentiable and  $\frac{d}{dx}(f(x)g(x)) = f(x)g'(x) + f'(x)g(x)$ .

As long as we are working our way through the combination of functions, we may as well tackle reciprocals and then quotients. Once again, we ask you to analyze a particular function in order to see the relationship among the derivative of the function, its function value, and the derivative of its reciprocal.

*Exercise 4.58.* Let  $f(x) = 3x$ . Using the definition of derivative, compute the value of  $\frac{d}{dx}(\frac{1}{f(x)})$  at the point  $x = 4$ . Once again as you do this exercise,



think through the definition of derivative to see how the derivative of the function, the value of the function, and the derivative of the reciprocal are related.

If you were successful with the previous exercise or if you remember the Reciprocal Rule from a calculus course, the following theorem will not be a surprise.

*Theorem 4.59 (Reciprocal Rule).* Let  $f(x)$  be a differentiable function with  $f(x_0) \neq 0$ . Then  $\frac{d}{dx}(\frac{1}{f(x)})_{x_0} = -\frac{f'(x_0)}{(f(x_0))^2}$ .

By combining the Reciprocal Rule and the Product Rule, we can formulate the quotient rule.

*Theorem 4.60 (Quotient Rule).* Let  $f(x)$  and  $g(x)$  be differentiable functions. Then  $\frac{d}{dx}(\frac{f(x)}{g(x)}) = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$  for every  $x$  for which  $g(x) \neq 0$ .

To actually take derivatives, the strategy is to individually take the derivatives of some basic functions using the definition of the derivative and then combine those results using rules of combination such as the sum, product, and quotient rules to determine the derivatives of more complicated functions. Let's now turn to trigonometric functions.

The trigonometric functions are differentiable, but again they present a challenge. Each of the basic trigonometric functions has its own difficulties, so let's just start with the sine and cosine. Prove the next two theorems using the definition of the derivative and the definitions of the sine and cosine. Draw a picture of the unit circle to see geometrically where every value in the definition of the derivative appears in the picture. Draw the picture big so that it is clear that a small segment of a circle looks like a straight line.

*Theorem\** 4.61. The trigonometric function  $\sin(x)$  is differentiable and

$$\frac{d}{dx}(\sin(x)) = \cos(x).$$

*Theorem\** 4.62. The trigonometric function  $\cos(x)$  is differentiable and

$$\frac{d}{dx}(\cos(x)) = -\sin(x).$$

We can now proceed to the other trigonometric functions by using the reciprocal and quotient rules.

*Exercise 4.63.* Derive the derivatives of the trigonometric functions  $\tan(x)$ ,  $\sec(x)$ ,  $\csc(x)$ , and  $\cot(x)$ .

One of the most potent methods for obtaining more complicated functions from simpler ones is to compose functions. Let's again see how the derivative of the composition of two functions is related to the derivatives of the two functions and their values.

*Exercise 4.64.* Let  $f(x) = x^2$  and  $g(x) = x^3$ . What is the derivative of the composition, that is,  $\frac{d}{dx}(g(f(x)))$ , at the point  $x = 4$ ? Of course, you could simply take the composition, which means to first square  $x$  to get  $x^2$  and then cube the result and realize that  $g(f(x)) = (x^2)^3 = x^6$  and then take its derivative. That is fine, particularly to check your thinking; however, for this exercise please think about the definition of the derivative. When you put in the  $h$  into the numerator of the definition, what happens? Try to follow the reasoning through, thinking about the derivatives of each of  $f(x)$  and  $g(x)$  and their values at relevant places as you do. The goal of this exercise is for you to see the relationships among the derivatives of each of the functions, the values of each of the functions, and the derivative of the composition.

If you were successful with the previous exercise or if you remember the Chain Rule from a calculus course, the following theorem will not be a surprise.

*Theorem 4.65 (Chain Rule).* Let  $f(x)$  and  $g(x)$  be differentiable functions from  $\mathbb{R}$  to  $\mathbb{R}$ . Then  $(g \circ f)(x) = g(f(x))$  is differentiable and  $\frac{d}{dx}(g(f(x))) = g'(f(x))f'(x)$ .

These theorems allow us to take derivatives of vast numbers of continuous functions such as  $\sin^3(x)\tan(x^2 + 3x + 4)$ .

Differentiable functions need not have domains equal to all of  $\mathbb{R}$ . For example, a function  $f : [0, 1] \rightarrow \mathbb{R}$  is defined only on the points in the interval  $[0, 1]$ . It is differentiable if and only if it is differentiable at each point of its domain. Since differentiable functions are continuous, we know that they attain their maximum and minimum values. If the maximal point does not occur at the endpoint, then they will have predictable derivatives.

*Theorem 4.66.* Let  $f : [0, 1] \rightarrow \mathbb{R}$  be a differentiable function. Suppose  $w \in (0, 1)$  is a real number such that for every  $x \in [0, 1]$ ,  $f(x) \leq f(w)$ . Then  $f'(w) = 0$ .

One theorem that captures the global implications of differentiability is the Mean Value Theorem, which implies that if Zeno went at a particular average velocity over an interval of time, then at some instant, his instantaneous velocity was that average velocity. This plausible statement can be couched in terms of derivatives.

*Theorem 4.67* (Mean Value Theorem). Let  $f : [a, b] \rightarrow \mathbb{R}$  be a differentiable function. Then for some real number  $c \in (a, b)$ ,  $f'(c) = \frac{f(b)-f(a)}{b-a}$ .

*Exercise 4.68.* Use the Mean Value Theorem to give a new proof that Zeno was speeding sometime between 3:00pm and 3:02pm.

## 4.7 Speedometer Movie and Position

This discussion of derivatives all emerged from solving the question of finding instantaneous velocity when we know the position at each moment. Let's return to moving cars to look at the reverse question, namely, finding the position if we know the instantaneous velocity at each moment.

Avant-garde movies strive for deep meaning, often with no action. These movies are incredibly boring and here we will describe some of the most boring. After their stint on the police force, Newton and Leibniz decided to turn their attention to film. They got in a car, turned the lens on the speedometer, and drove forward on a straight road for an hour. The movie was not edited and presented only the speedometer dial with the needle sometimes moving slowly, sometimes fixed for minutes on end. None of the road could be seen and the action was unrelieved by a glimpse at the odometer. The movie was time-stamped at each moment, so the viewer could see how much of life would be wasted before the merciful conclusion of this 'drama'. Newton and Leibniz made several of these hour long movies; however, few people went back to see the sequels.

Since viewers were terminally bored with these movies, Newton and Leibniz decided to pose a question to give their audience something to do. They asked when the movie began, "How far did the car go during this hour?"

This question turned the movie from a sleeper to a riveting challenge that changed the world.

*Exercise 4.69.* Here are some descriptions of the speedometer movies. For each one, figure out how far the car went and develop a method that would work for any such movie.

1. This movie was the most boring of all. For the entire movie, the speedometer read 30 mph.
2. This movie had only one change. For the first half hour, the speedometer read 30 MPH, and then instantly changed to read 60 MPH for the second half hour.

3. In this movie, the speedometer started at 0 MPH and gradually and uniformly increased by 1 MPH each minute to read 60 MPH at the end of the hour.
4. In this movie, the speedometer's reading was always  $t^2$  where  $t$  was the number of minutes into the movie. This car was really moving by the end of the hour.
5. In this movie, the movie was changing speeds all the time. What strategy could you devise to pin down the distance traveled within 1 mile? ...within 0.1 miles? ...within .001 miles? ...to pin down the distance exactly?

In answering the previous exercise, you have defined the definite integral. In the following definition, think of the function  $f(x)$  as telling the speedometer reading at each time  $x$ .

*Definition.* Let  $f(x)$  be a continuous function on the interval  $[a, b]$ . Then the **definite integral** of  $f(x)$  from  $a$  to  $b$  is a limit of a sequence of approximating sums of products where the  $n$ th approximation is obtained by dividing the interval  $[a, b]$  into  $n$  equal subintervals:  $[a_0 = a, a_1]$ ,  $[a_1, a_2]$ ,  $[a_2, a_3], \dots, [a_{n-2}, a_{n-1}], [a_{n-1}, a_n = b]$ , then for each subinterval multiplying its width (which is always  $\frac{b-a}{n}$ ) by the value of the function at its end-point (that product would give approximately the distance traveled during that small interval of time) and then adding all those products up. For every choice of  $n$  number of intervals we have an approximation of the integral, so the integral is the limit as we choose increasingly large  $n$ , which produces increasing small subintervals. So in symbols:  $\int_a^b f(x) dx = \lim_{n \rightarrow \infty} \sum_{k=0}^{n-1} f(a + k \frac{b-a}{n}) \frac{b-a}{n}$ .

Leibniz is again responsible for the notation for the interval. Notice that every feature of the notation refers to its definition. The long s shape stands for 'sum', the limits of integration tell us where the  $x$  is varying between, the ' $dx$ ' is the small width and it is next to the  $f(x)$ , so  $f(x)dx$  suggests the distance traveled in the small ' $dx$ ' interval of time. So adding up those small contributions to the distance traveled gives the total distance traveled.

## 4.8 Fundamental Theorem of Calculus

Since the derivative and the integral really involved the same car moving down the road, there is a clear and natural connection between the two concepts of the derivative and the integral. Namely, there are two ways to

look at how far the car traveling along a straight road has traveled. On the one hand, see where the car was at the end and subtract where it was at the beginning in order to compute the net change over that interval of time. The other way is to do the integral procedure. Since both methods yield the same result of the net change in the position of the car, those two method must produce the same answer. But notice that if we have a position function  $p(t)$  that is telling us the position of the car at every time  $t$  from some time  $a$  to time  $b$ , then  $p'(t)$  is telling us what the speedometer will be reading at each moment. So we can see that the integral of  $p'(t)$  from time  $a$  to time  $b$  will give the same answer as the difference in the ending position  $p(b)$  minus the starting position  $p(a)$ . This insight is the most important insight in calculus and therefore has the exalted title of The Fundamental Theorem of Calculus.

*Theorem 4.70* (Fundamental Theorem of Calculus). Let  $F(x)$  be a function on the interval  $[a, b]$  with derivative  $F'(x)$ . Then  $\int_a^b F'(x)dx = F(b) - F(a)$ .

The definition of the definite integral tells us that the value of the integral is something meaningful that we want to know, such as the distance the car has traveled. The Fundamental Theorem of Calculus tells us that to find the value of a definite integral, all we need to do is to find an antiderivative, plug into two values, and subtract. So the Fundamental Theorem of Calculus is the reason that anti-derivatives are so closely linked with integrals. In fact, we soon start saying 'integral' when we mean 'anti-derivative'.

The end

---

Integrals, Mean Value for Integrals

Derivative theorems and formulas, integral formulas



## Chapter 5

# Topology - Math from Math

Our exploration of graph theory arose from looking at the real-world Koenigsberg Bridge Problem. The ideas of group theory emerged from isolating commonalities among several familiar examples such as motions of a triangular block. The epsilon-delta definitions of calculus came from analyzing the motion of a car. All of those wonderful mathematical adventures began with experiences from the real world. The world is the ultimate source of mathematical inspiration. However, once we have created mathematical concepts, those concepts become part of our world. They have a life of their own. They have all the richness and subtlety that our explorations have suggested—and much more. Those mathematical worlds become in a sense as much a part of the real world of our minds as is the real-world that we perceive more directly. In reality, all of our experience of the world occurs in our minds, so in some real sense, there is little actual distinction between the ‘real world’ and the ‘abstract’ world of mathematical ideas that we have created. Therefore, it should come as no surprise that one of the most fruitful sources for mathematical inspiration is mathematics itself. We can investigate mathematics that we know and from it seek essential ingredients and commonalities, which then give rise to concepts that we can explore. Topology is a mathematical subject that arose from reflecting on the world of mathematics itself.

At about the turn of the 20th century, that is, about the year 1900, mathematicians undertook to put all of mathematics on a firmer foundation. The ancient model for mathematical rigor was Euclid’s “Elements”. Euclid presented geometry and other mathematics by beginning with a set of definitions and axioms from which all the theorems followed logically. So at the turn of the 20th century, mathematicians sought to duplicate Euclid’s

strategy by seeking a set of definitions and axioms from which mathematical ideas such as continuity and other analytic ideas would follow. At that time, mathematicians realized that we have to start someplace with undefined terms.

The primitive undefined terms for this program of axiomatizing mathematics were point and set. And the goal was to begin with points and sets and capture the reason for the richness of calculus using those primitive ideas. Isolating set-theoretic principles that captured notions of convergence and continuity was the triumph of topology.

## 5.1 Closeness

The calculus concepts of continuity, derivative, and integral hinge on ideas of convergence. Convergence depends on a concept of real numbers getting closer and closer. When we define the derivative, we take average speeds or average slopes over increasingly smaller intervals and find that those averages converge to a single value. The approximations get closer and closer to a single value.

Of course, it is natural to think of the term ‘close’ and the concept of nearness as an expression of some underlying sense of distance. But let’s step back and consider the essential use that we are making of the idea of closeness. When we say, “Consider all the real numbers that are less than 0.2 from 5,” we are actually describing a set of points, in this case, the interval  $(4.8, 5.2)$ . So from one point of view, the set of all points less than 0.2 from 5 is describing an example of a set in which we have a special interest. The points less than 0.1 from 6 is another set of interest. So sets that are defined as those of less than a certain distance from a given point are special sets. These sets are subsets of the underlying total space, in this case, the real line. Some subsets of the reals are not pertinent to our discussions about convergence and continuity. For example, the set of rational numbers is not a set that captures our concept of closeness.

Now let’s try to describe convergence and continuity without using distances. In general terms, a continuous function is one in which close points in the domain are close in the range. We can take one further step toward the precise definition of continuity by saying that for a continuous function, given any distance in the range around  $f(x)$ , there is a distance in the domain such that all points less than that distance in the domain from  $x$  get mapped to within the given distance in the range around  $f(x)$ . Let’s rephrase that definition without using the idea of distance. We could say, a function  $f$  is



continuous at  $x$  if given any special set around  $f(x)$  in the range, there is a special set in the domain around  $x$  such that all the points in the domain special set get mapped into the range special set.

We could phrase the concept of convergence in that same general manner. We would say that an infinite list of points converges to a limit point means that for any special set around the limit point, all but a finite number of points in our list of points lie in that special set.

Of course, if all we meant by ‘special set’ was  $\epsilon$  neighborhoods, we would not have accomplished our goal of generalizing the ideas of convergence or continuity. We would simply have rephrased them. So now we break free of the concept of distance and think far more abstractly. The concept of convergence basically is saying that in some set  $X$  we have designated a special collection of subsets that we will use to define the meaning of convergence. A set of points converges to a limit point if for any one of our special sets around the limit point, all but finitely many of our set of points lies inside that special set.

Thinking about continuity in the abstract, we have two sets  $X$  and  $Y$ , the domain  $X$  and the range  $Y$ . There are some specially designated subsets of  $X$  and specially designated subsets of  $Y$ . Then we want to say that a function  $f : X \rightarrow Y$  is continuous if for any point  $x \in X$  and any specially designated subset  $U$  of  $Y$  containing  $f(x)$ , there is a specially designated subset  $V$  of  $X$  containing  $x$  such that  $f(V) \subseteq U$ .

These specially designated subsets are not just  $\epsilon$ -neighborhoods, but they are generalizations of that idea.

## 5.2 Definition of a Topology

So far we have described a format for the phraseology of concepts like convergence and continuity; however, we have not defined the conditions that would make a reasonable collection of distinguished subsets. Again let’s return to our generative example of the real numbers with our usual sense of distance and see whether there are some conditions on the distinguished subsets that we want. We want our collection of distinguished subsets to reflect the essential features that make convergence and continuity work.

In describing convergence and continuity, we found ourselves selecting a challenge set around a point and then concluding that some points were in that challenge set. In the case of convergence, all except a finite number of points were in the challenge set. So if we have two challenge sets containing the same point, then all except a finite number of points must be in the

intersection of the two challenge sets. This suggests that the intersection of two distinguished sets should also be a distinguished set. We can reach the same intuition when we think about continuity.

Let's think about a function  $f : X \rightarrow Y$ . Recall that our intuition told us that  $f$  is continuous if for any point  $x \in X$  and any specially designated subset  $U$  of  $Y$  containing  $f(x)$ , there is a specially designated subset  $V$  of  $X$  containing  $x$  such that  $f(V) \subseteq U$ . However, we might want to think more globally and realize that many points may all go to the same point  $y$  in  $Y$ . So if  $f$  is continuous and  $y$  is a point in  $Y$  and  $V$  is a specially designated set containing  $y$ , then for each such point  $x$  for which  $f(x)$  goes to  $y$ , there is a designated set  $U_x$  in  $X$  with  $x$  in  $U_x$  such that  $f(U_x) \subseteq V$ . That means the union of all such designated subsets  $U_x$  must also have the property that that union's image is in  $V$ . This property suggests that arbitrary unions of distinguished sets should be declared to be distinguished.

The two properties that finite intersections of distinguished subsets and arbitrary unions of distinguished subsets are distinguished form the essential ingredients in formulating the definition of a topology. The definition is completed by making the whole space  $X$  be distinguished, which is equivalent to saying that each point of  $X$  is in some one of the distinguished subsets. Finally, for technical reasons it is convenient to include the empty set as distinguished. We are now ready for the formal definition of a topology.

*Definition.* Let  $X$  be a set and  $\mathcal{T}$  a collection of subsets of  $X$ . Then  $\mathcal{T}$  is a *topology* on  $X$  if and only if

1.  $\emptyset \in \mathcal{T}$ ,
2.  $X \in \mathcal{T}$ ,
3. if  $U \in \mathcal{T}$  and  $V \in \mathcal{T}$ , then  $U \cap V \in \mathcal{T}$ , and
4. if  $\{U_\alpha\}_{\alpha \in \lambda}$  is any collection of sets each in  $\mathcal{T}$ , then  $\cup_{\alpha \in \lambda} U_\alpha \in \mathcal{T}$ .

A *topological space* is a pair  $(X, \mathcal{T})$  where  $X$  is a set and  $\mathcal{T}$  is a topology for  $X$ . If  $(X, \mathcal{T})$  is a topological space, then  $U \subseteq X$  is called an *open set* in  $(X, \mathcal{T})$  if and only if  $U \in \mathcal{T}$ . We sometimes say, ' $X$  is a topological space.' When we say that, we mean that there is a topology  $\mathcal{T}$  on  $X$  that is implicit.

And now we must translate our notion of "closeness".

*Definition.* Let  $(X, \mathcal{T})$  be a topological space. Let  $A$  be a subset of  $X$ , written  $A \subset X$ ; and let  $p$  a point in  $X$ , written  $p \in X$ . Then  $p$  is a *limit point* of  $A$  if, for every  $U \in \mathcal{T}$  containing  $p$ ,  $(U \setminus \{p\}) \cap A \neq \emptyset$ . Notice that  $p$  may or may not be in  $A$ .

In other words,  $p$  is a limit point of  $A$  if *all* open sets containing  $p$  intersect  $A$  at *some point other than* (possibly)  $p$  itself. Or we might say the point  $p$  is a limit point of  $A$  if it cannot be separated from  $A$  by an open set. A point  $p$  is close to a set  $A$  if  $p$  is a limit point of  $A$  in exactly the same way that the border of Texas is close to Texas because every map containing the border also contains points inside Texas. Thus, the concept of a limit point gives us a way of capturing the idea of a point “being arbitrarily close” to a set *without* using the concept of distance. Instead we use the idea of open sets in a topology.

Notice that the idea of being a limit point is somewhat different from the idea of convergence in that being a limit point does not insist that all the points except for finitely many lie in the open set, but merely that every open set around the limit point contains at least one point of the other set.

The first step towards understanding such an abstract definition is to generate a few examples and carefully check, using the technical definition, that they have all of the desired properties.

The first example is an extreme. It is extreme in its having the fewest possible open sets, namely only the required whole space and empty set.

*Example 1* (indiscrete topology). For a set  $X$ ,  $\mathcal{T} = \{\emptyset, X\}$  is called the *indiscrete topology* on  $X$ . So  $(X, \{\emptyset, X\})$  is an indiscrete topological space.

What would it mean to check that  $\mathcal{T} = \{\emptyset, X\}$  is actually a topology on some set  $X$ ? Well, going back to the definition of a topology, certainly  $X$  is a set and  $\mathcal{T}$  a collection of subsets of  $X$ , so there are four conditions that we must check. The first two properties of a topology are true by inspection, but the other two involve some work. This amounts to understanding what the symbols  $\cap$  (intersection) and  $\cup$  (union) really mean.

*Exercise 5.1.* Let  $X$  be a set and  $\mathcal{T} = \{\emptyset, X\}$ . Check that  $\mathcal{T}$  is a topology on  $X$ .

Our favorite set in this chapter will be  $\mathbb{R}$ , the real numbers, sometimes called the real line. So when you see a set  $X$  in a definition, imagine that set to be  $\mathbb{R}$  for your first example. Since we will be writing lots of subsets of the real line, we should fix some notation. We will write  $(a, b)$  for the set  $\{r \in \mathbb{R} \mid a < r < b\}$  and  $[a, b]$  for the set  $\{r \in \mathbb{R} \mid a \leq r \leq b\}$ . Similarly  $(a, b] = \{r \in \mathbb{R} \mid a < r \leq b\}$  and  $[a, b) = \{r \in \mathbb{R} \mid a \leq r < b\}$ . We will call each of these sets an interval with endpoints  $a$  and  $b$ . To tell them apart in speech, we might call  $[a, b)$  “the interval from  $a$  to  $b$  including  $a$ ”. Notice that if  $b < a$ , then any of these sets is empty, so we will consider the empty set to be an interval. Also, we allow  $\pm\infty$  as the end of an interval; but

$\pm\infty$  is never in the interval, so it always appears with a curved bracket, eg  $(-\infty, 2]$ .

The set  $\mathbb{R}$  will be our favorite set; however, we will use  $\mathbb{R}$  in unfamiliar ways. Specifically, the open intervals that we are accustomed to may or may not be open sets in the various different topologies that we will consider for  $\mathbb{R}$ . Let's investigate  $\mathbb{R}$  with the indiscrete topology, for example.

*Exercise 5.2.* Consider  $\mathbb{R}$  with the indiscrete topology  $\mathcal{T} = \{\emptyset, \mathbb{R}\}$ . Is  $p = 1$  a limit point of the set  $A = [0, 1)$ ? What about if  $p = 0$  or  $p = -1$ ?

We say that  $(X, \mathcal{T})$  and  $(X, \mathcal{T}')$  are *different* topological spaces if  $\mathcal{T} \neq \mathcal{T}'$ , even though the underlying set  $X$  is the same. Keep in mind that open sets  $U$  are *elements* ( $\in$ ) of the topology  $\mathcal{T}$  and *subsets* ( $\subset$  or  $\subseteq$ ) of the space  $X$ , and we usually pick capitol letters to represent them. Elements of  $X$ , on the other hand, are what we call the points of the space  $X$ , usually denoted by lowercase letters like  $p$ .

What is the biggest collection of open sets that any topology could have? The answer is that we could make every subset of  $X$  a member of the topology.

*Example 2* (discrete topology). For a set  $X$ , let  $2^X$  be the set of all subsets of  $X$ . Then  $\mathcal{T} = 2^X$  is called the *discrete topology* on  $X$ . The space  $(X, 2^X)$  is called a *discrete topological space*. Note the spelling: *discrete* topology, not *discreet* topology!

*Exercise 5.3.* Let  $X$  be any set and  $\mathcal{T} = 2^X$ . Show that  $(X, \mathcal{T})$  is a topological space.

*Exercise 5.4.* Consider  $\mathbb{R}$  with the discrete topology  $\mathcal{T} = 2^{\mathbb{R}}$ . Is  $p = 1$  a limit point of the set  $A = [0, 1)$ ? What about if  $p = 0$  or  $p = -1$ ?

Notice that the discrete topology has the maximum possible collection of open sets that any topology can have while the indiscrete topology has the minimum possible collection of open sets. Now that we have investigated the definition a little by checking the extreme cases, let's consider a familiar example and describe a natural topology for it. Imagine our set to be a plane, in particular  $\mathbb{R}^2$ , the real plane. A subset will be declared to be an open set, that is, a member of the topology, if around every point in the set, there is a little rectangle, a zoomed in map, entirely contained within the set.

*Example 3* (standard topology on  $\mathbb{R}^2$ ). The *standard topology*  $\mathcal{T}_{\text{std}}$  for  $\mathbb{R}^2$  is defined as follows: a subset  $U$  of  $\mathbb{R}^2$  belongs to  $\mathcal{T}_{\text{std}}$  if and only if for each point  $p$  of  $U$  there is a rectangle,  $(a, b) \times (c, d)$ , containing  $p$ , such that the rectangle is entirely contained in  $U$ .

When we talk about rectangles, we do not include their borders, only the interiors. And we are only considering rectangles whose sides are parallel to the axes.

*Exercise 5.5.* Show that  $(\mathbb{R}^2, \mathcal{T}_{\text{std}})$  is a topological space.

*Exercise 5.6.* Consider  $\mathbb{R}^2$  with the standard topology  $\mathcal{T} = \mathcal{T}_{\text{std}}$ . Is  $p = (1, 1)$  a limit point of the unit square (excluding its boundary)  $A = \{(x, y) \mid 0 < x < 1, 0 < y < 1\}$ ? What about if  $p = (1, \frac{1}{2})$  or  $p = (0, -1)$ ?

There is nothing special about the 2 (in the  $\mathbb{R}^2$ ) in this last example; there is an analogous definition for any  $\mathbb{R}^n$ . But what will play the role of the rectangle? Well, if the four corners of a rectangle are the points  $Q = (q_1, q_2)$  (lower left corner),  $R = (r_1, r_2)$  (lower right corner),  $S = (s_1, s_2)$  (upper right corner), and  $T = (t_1, t_2)$  (upper left corner), then the rectangle is the points  $(x, y)$  such that  $q_1 < x < r_1$  and  $q_2 < y < t_2$ . Then the 1-dimensional analogue should be all points  $x$  such that  $a < x < b$ , namely an interval. Similarly, a the 3-dimensional version is a box (for example a cube). We will call these objects  $n$ -boxes, where  $n$  is the dimension.

*Definition.* An  $n$ -box is a set of the form  $\{(x_1, x_2, \dots, x_n) \mid a_i < x_i < b_i\}$ , where  $a_i$  and  $b_i$  are real numbers. Notice that if some  $b_i \leq a_i$ , then this set is empty, so we consider the empty set to be an  $n$ -box for every  $n$ .

*Example 4* (standard topology on  $\mathbb{R}^n$ ). The *standard topology*  $\mathcal{T}_{\text{std}}$  for  $\mathbb{R}^n$  is defined as follows: a subset  $U$  of  $\mathbb{R}^n$  belongs to  $\mathcal{T}_{\text{std}}$  if and only if for each point  $p$  of  $U$  there is an  $n$ -box containing  $p$ , such that the  $n$ -box is entirely contained in  $U$ .

In particular, a set  $U$  is open in the standard topology on  $\mathbb{R}$ , the real line, if for every  $p \in U$ , there exist numbers  $a_p$  and  $b_p$  such that  $p \in (a_p, b_p) \subset U$ .

*Theorem 5.7.* A set  $U \subset \mathbb{R}^n$  is open in the standard topology if and only if it can be written as a union of  $n$ -boxes.

There are a few other common topologies that are interesting to play with.

*Example 5* (finite complement or co-finite topology). For any set  $X$ , the *finite complement (or co-finite) topology* for  $X$  is described as follows: a subset  $U$  of  $X$  is open if and only if  $U = \emptyset$  or  $X - U$  is finite.

*Exercise 5.8.* Verify that all the examples given above are indeed topologies; in other words, that they satisfy all four conditions needed to be a topology.

*Exercise 5.9.* 1. Describe some of the open sets you get if  $\mathbb{R}$  is endowed with the topologies described above (standard, discrete, indiscrete,

and co-finite). Specifically, identify sets that demonstrate the differences among these topologies, that is, find sets that are open in some topologies but not in others.

2. For each of the topologies, determine if the interval  $(0, 1) \in \mathbb{R}$  is an open set in that topology.

Intuitively, points inside a set should be close to that set, but that's not quite how our definition works.

*Definition* (isolated point). Let  $(X, \mathcal{T})$  be a topological space,  $A$  be a subset of  $X$ , and  $p$  be a point in  $X$ . If  $p \in A$  but  $p$  is not a limit point of  $A$ , then  $p$  is an *isolated point* of  $A$ .

If  $p$  is an isolated point of  $A$ , then there is an open set  $U$  such that  $U \cap A = \{p\}$ .

*Theorem 5.10.* Suppose  $A$  is a subset of  $X$  and  $p \notin A$  is a point in a topological space  $(X, \mathcal{T})$ . Then  $p$  is *not* a limit point of  $A$  if and only if there exists an open set  $U$  with  $p \in U$  and  $U \cap A = \emptyset$ . (Careful, this looks a lot like the definition, but it's not.)

*Exercise 5.11.* Give examples of a set  $A$  in a topological space and

1. a limit point of  $A$  that is an element of  $A$ ;
2. a limit point of  $A$  that is not an element of  $A$ ;
3. an isolated point of  $A$ ;
4. a point not in  $A$  that is not a limit point of  $A$ .

We finish this section by investigating why we restrict ourselves to the intersection of two open sets in the definition of a topology.

*Exercise 5.12.* Give an example of a topological space and a collection of open sets in that topological space to show that the *infinite* intersection of open sets need not be open.

*Theorem 5.13.* Let  $\{U_i\}_{i=1}^n$  be a *finite* collection of open sets in a topological space  $(X, \mathcal{T})$ . Then  $\bigcap_{i=1}^n U_i$  is open.

### 5.3 Closed Sets

We've been talking about sets and their limit points. What properties does a set have if it already contains all of its limit points? In other words what kinds of sets already contain all points that are near them?

*Definition* (closure of a set). Let  $(X, \mathcal{T})$  be a topological space, and  $A \subseteq X$ . Then the *closure* of  $A$ , denoted  $\overline{A}$  or  $\text{Cl}(A)$ , is  $A$  together with all of its limit points.

*Definition* (closed set). Let  $(X, \mathcal{T})$  be a topological space and  $A \subseteq X$ .  $A$  is *closed* if and only if  $\text{Cl}(A) = A$ , in other words, if  $A$  contains all its limit points.

The reason that the next theorem is not obvious is that we might think that when we add limit points of a set the new augmented set may then have yet more limit points; however, that possibility cannot occur, as this next theorem proves.

*Theorem 5.14.* For any topological space  $(X, \mathcal{T})$  and  $A \subseteq X$ ,  $\overline{A}$  is closed, that is, for any set  $A$  in a topological space,  $\text{Cl}(\overline{A}) = \overline{A}$ .

Remember that we think of the limit points of a set  $A$  as the points that are close to  $A$ . So, if a point  $p$  is close to a set  $A$ , and  $A$  is contained in a set  $B$ , then  $p$  ought to be close to  $B$  as well. Also, the only way to be close to the union of two sets is to be close to at least one of them. This argument is, of course, not a proof.

*Theorem 5.15.* Let  $A, B$  be subsets of a topological space. Then

1.  $A \subseteq B \Rightarrow \overline{A} \subseteq \overline{B}$ ; and
2.  $\overline{A \cup B} = \overline{A} \cup \overline{B}$ .

A basic relationship between open sets and closed sets in a topological space is that they are complements of each other. If you are not familiar with the phrase “complement of  $A$  (in  $X$ )”, it simply means the set that contains everything outside of  $A$ , sometimes written  $X \setminus A$ .

*Theorem 5.16.* Let  $(X, \mathcal{T})$  be a topological space. Then the set  $A$  is closed if and only if  $X - A$  is open.

When we remove a closed set from an open set, we are left with an open set, and when we remove an open set from a closed set, the result is still closed.

*Theorem 5.17.* Let  $(X, \mathcal{T})$  be a topological space, and let  $U$  be an open set and  $A$  be a closed subset of  $X$ . Then the set  $U - A$  is open and the set  $A - U$  is closed.

Notice that most subsets of a topological space  $X$  are typically neither open nor closed. It is wrong to say, “Well, it’s not open, so it must be closed.”

*Exercise 5.18.* Show that the subset consisting of the rational numbers is neither an open nor a closed set in the standard topology on  $\mathbb{R}$ .

The properties of a topological space can be captured by focusing on closed sets instead of open sets. From that perspective, the four defining properties of a topological space are captured in the following theorem about closed sets where only finite unions are allowed, but arbitrary intersections are permitted.

*Theorem 5.19.* Let  $(X, \mathcal{T})$  be a topological space:

1.  $\emptyset$  is closed.
2.  $X$  is closed.
3. The union of finitely many closed sets is closed.
4. Let  $\{A_\alpha\}_{\alpha \in \lambda}$  be a collection of closed subsets in  $(X, \mathcal{T})$ . Then  $\bigcap_{\alpha \in \lambda} A_\alpha$  is closed.

*Exercise 5.20.* Give an example to show that the union of infinitely many closed sets in a topological space may be a set that is not closed.

*Exercise 5.21.* Give examples of topological spaces and sets in them that:

1. are closed, but not open;
2. are open, but not closed;
3. are both open and closed;
4. are neither open nor closed.

*Exercise 5.22.* State whether each of the following sets are open, closed, both or neither.

1. In  $\mathbb{Z}$  with the finite complement topology:  $\{0, 1, 2\}$ ,  $\{\text{prime numbers}\}$ ,  $\{n : |n| \geq 10\}$ .
2. In  $\mathbb{R}$  with the standard topology:  $A = (0, 1)$ ,  $B = (0, 1]$ ,  $C = [0, 1]$ ,  $D = \{0, 1\}$ ,  $E = \{\frac{1}{n} | n \in \mathbb{N}\}$ .
3. In  $\mathbb{R}^2$  with the standard topology:  $C = \{(x, y) | x^2 + y^2 = 1\}$ ,  $D = \{(x, y) | x^2 + y^2 > 1\}$ ,  $\Omega = \{(x, y) | x^2 + y^2 \geq 1\}$ ,



4. Which subsets are closed in a set  $X$  with the discrete topology? indiscrete topology?

*Theorem 5.23.* For any set  $A$  in a topological space  $X$ , the closure of  $A$  equals the intersection of all closed sets containing  $A$ , that is,

$$\text{Cl}(A) = \bigcap_{A \subseteq C, C \in \mathcal{C}} C$$

where  $\mathcal{C}$  is the collection of all closed sets in  $X$ .

Informally, we can say  $\overline{A}$  is the “smallest” closed set that contains  $A$ .

*Exercise 5.24.* Pick several different subsets of  $\mathbb{R}$ , and find their closure in:

1. the discrete topology;
2. the indiscrete topology;
3. the finite complement topology;
4. the standard topology.

*Exercise 5.25.* In  $\mathbb{R}^2$  with the standard topology, describe the limit points and closure of the following two sets:

1. The topologist’s sine curve:

$$S = \left\{ \left( x, \sin \left( \frac{1}{x} \right) \right) \mid x \in (0, 1) \right\}$$

2. The topologist’s comb:

$$C = \{(x, 0) \mid x \in [0, 1]\} \cup \bigcup_{n=1}^{\infty} \left\{ \left( \frac{1}{n}, y \right) \mid y \in [0, 1] \right\}$$

The following exercise is difficult.

*Exercise 5.26.* In the standard topology on  $\mathbb{R}$ , describe a non-empty subset  $C$  of the closed unit interval  $[0, 1]$  that is closed, contains no non-empty open interval, and where no point of  $C$  is an isolated point.

## 5.4 Subspaces

If  $(X, \mathcal{T}_X)$  is a topological space and  $Y$  is a subset of  $X$ , then there is a natural topology on  $Y$  induced by  $\mathcal{T}$ :

*Definition* (subspace topology). Let  $(X, \mathcal{T}_X)$  be a topological space. For  $Y \subseteq X$ , the collection

$$\mathcal{T}_Y = \{U \mid U = V \cap Y \text{ for some } V \in \mathcal{T}_X\}$$

is a topology for  $Y$ , called the *subspace topology*. The space  $(Y, \mathcal{T}_Y)$  is called a (topological) *subspace* of  $X$ .

*Theorem 5.27.* The subspace topology defined above,  $\mathcal{T}_Y$ , is in fact a topology.

*Exercise 5.28.* In  $Y = (0, 1)$ , as a subspace of  $\mathbb{R}_{\text{std}}$ , is  $(\frac{1}{2}, 1)$  closed, open, both, or neither? What about  $[\frac{1}{2}, 1)$ ?

Now that there are multiple topologies floating around in the same question, we will have to be careful about our language. When we say that a set  $U$  is open, we must state with respect to which topology. For example, we could say that  $U$  is open in  $X$  or  $U \in \mathcal{T}_X$ . Similarly the closure of a set depends on which topology you're using, so we must specify. For example, we could say the closure of  $A$  in  $X$  or  $\text{Cl}_X(A)$ .

*Exercise 5.29.* Consider a subspace  $(Y, \mathcal{T}_Y) \subseteq (X, \mathcal{T}_X)$ . Is every subset  $U \subseteq Y$  that is open with respect to the subspace topology also open in  $(X, \mathcal{T}_X)$ ?

If so, prove it. If not, give an example of a set  $Y$  and a subset  $U$  that is open in  $Y$  but not open in  $X$ .

*Theorem 5.30.* Let  $(Y, \mathcal{T}_Y)$  be a subspace of  $(X, \mathcal{T})$ . A subset  $A$  is closed in  $(Y, \mathcal{T}_Y)$  if and only if there is a set  $B \subset X$ , closed in  $X$ , such that  $A = Y \cap B$ .

*Theorem 5.31.* Let  $(Y, \mathcal{T}_Y)$  be a subspace of  $(X, \mathcal{T}_X)$ . A subset  $A \subset Y$  is closed in  $(Y, \mathcal{T}_Y)$  if and only if  $\text{Cl}_X(A) \cap Y = A$ .

A stronger version of the previous theorem would state that for a subspace  $(Y, \mathcal{T}_Y) \subset (X, \mathcal{T}_X)$  and a set  $A \subset Y$ ,  $\text{Cl}_Y(A) = \text{Cl}_X(A) \cap Y$ .

## 5.5 Bases

Because arbitrary unions of open sets are open, a topological space can have extremely complicated open sets. It is often convenient to describe a (simpler) subcollection of open sets that *generate* all open sets in a given topology. Generating the group in group theory context was done using

repeated applications of the binary operation using a subset of the group elements. Generating a topology in this topological context will be done by taking arbitrary unions of a subcollection of the open sets in the topology. So instead of having to specifically describe all of the open sets in a topological space  $(X, \mathcal{T})$ , we can more conveniently specify a subcollection, called a *basis* for the topology  $\mathcal{T}$ . Recall, for instance, that in order to define the open sets in the standard topology in  $\mathbb{R}^2$  (respectively,  $\mathbb{R}^n$ ) we used the concept of rectangles (respectively,  $n$ -boxes). We called a set  $U$  an open set if, for each point  $p$  in  $U$ , we could find a rectangle (respectively,  $n$ -box) containing  $p$  contained in  $U$ . Thus, we could think of open sets as being made by taking arbitrary unions of these simpler open sets, the boxes.

*Definition* (basis of a topology). Let  $\mathcal{T}$  be a topology on a set  $X$  and let  $\mathcal{B} \subseteq \mathcal{T}$ . Then  $\mathcal{B}$  is a *basis* for the topology  $\mathcal{T}$  if and only if every element of  $\mathcal{T}$  is the union of elements in  $\mathcal{B}$ . If  $B \in \mathcal{B}$ , we say  $B$  is a *basis element* or *basic open set*. Note that  $B$  is an *element* of the basis, but a *subset* of the space. If  $\mathcal{B}$  is a basis for a topology  $\mathcal{T}$ , then we say that  $\mathcal{T}$  is *generated* by  $\mathcal{B}$  and an element  $U \in \mathcal{T}$  is generated by  $\mathcal{B}$  if it can be written as a union of elements in  $\mathcal{B}$ .

*Theorem 5.32.* Let  $(X, \mathcal{T})$  be a topological space and  $\mathcal{B}$  be a collection of subsets of  $X$ . Then  $\mathcal{B}$  is a basis for  $\mathcal{T}$  if and only if

1.  $\mathcal{B} \subseteq \mathcal{T}$ ,
2.  $\emptyset \in \mathcal{B}$ ,
3. for each set  $U$  in  $\mathcal{T}$  and point  $p$  in  $U$  there is a set  $B$  in  $\mathcal{B}$  such that  $p \in B \subseteq U$ .

Just as any group has a set of generators, namely the collection of *all* elements, every topology  $\mathcal{T}$  has a basis, namely  $\mathcal{T}$ . But this basis isn't a simplification. Fortunately, the topologies we've discussed have interesting bases.

*Theorem 5.33.* Let  $\mathcal{B}_1 = \{(a, b) \subseteq \mathbb{R} \mid a, b \in \mathbb{Q}\}$ , then  $\mathcal{B}_1$  is a basis for  $\mathbb{R}_{std}$ . Let  $\mathcal{B}_2 = \{(a, b) \cup (c, d) \subseteq \mathbb{R} \mid a, b, c, d \text{ are distinct irrational numbers}\}$ , then  $\mathcal{B}_2$  is also a basis for  $\mathbb{R}_{std}$ , the *standard topology* on  $\mathbb{R}$ .

*Exercise 5.34.* Describe a basis for the discrete topology on  $\mathbb{R}$ . What is the smallest possible collection that forms a basis?

*Theorem 5.35.* Let  $(X, \mathcal{T}_X)$  be a topological space, and  $(Y, \mathcal{T}_Y)$  be a subspace. If  $\mathcal{B}$  is a basis for  $\mathcal{T}_X$ , then  $\mathcal{B}_Y = \{B \cap Y \mid B \in \mathcal{B}\}$  is a basis for  $\mathcal{T}_Y$ .

Sometimes it is convenient to describe a topological space by describing a basis and then checking that the collection of all unions of those subsets does in fact satisfy the definition of a topology. The Lower Limit Topology is best described in this way.

*Example 6* (lower limit topology). We can define another topology on  $\mathbb{R}$ , called the *lower limit topology*, generated by a basis consisting of all sets of the form  $[a, b) = \{x \in \mathbb{R} \mid a \leq x < b\}$ . Denote this space by  $\mathbb{R}_{LL}$ . The real line with the lower limit topology is sometimes called the *Sorgenfrey line* or  $\mathbb{R}^1(\text{bad})$ .

*Theorem 5.36.* Show that  $\mathbb{R}_{LL}$  is a topological space.

- Exercise 5.37.*
1. Show that every open set in the standard topology on  $\mathbb{R}$  is open in  $\mathbb{R}_{LL}$ .
  2. Describe some sets that are open in  $\mathbb{R}_{LL}$  that are not open in the standard topology on  $\mathbb{R}$ .
  3. Describe an infinite set in  $[0, 1]$  that has no limit point in the  $\mathbb{R}_{LL}$  topology.

## 5.6 Product Topologies

The Cartesian product of two topological spaces has a natural topology derived from the topologies on each one of the component spaces. Let's recall the definition of the Cartesian product.

*Definition* (product). Let  $X$  and  $Y$  be two sets. The *product*  $X \times Y$ , or *Cartesian product*, is the set of ordered pairs  $(x, y)$  where  $x \in X$  and  $y \in Y$ .

*Exercise 5.38.* Give an example of  $(X, \mathcal{T}_X)$  and  $(Y, \mathcal{T}_Y)$ , topological spaces, such that  $\mathcal{T} = \{U \times V \mid U \in \mathcal{T}_X, V \in \mathcal{T}_Y\}$  is *not* a topology on  $X \times Y$ .

*Definition* (product topology). Suppose  $X$  and  $Y$  are topological spaces. The *product topology*,  $\mathcal{T}_{X \times Y}$ , on the product  $X \times Y$  is the topology whose *basis* is all sets of the form  $U \times V$  where  $U$  is an open set in  $X$  and  $V$  is an open set in  $Y$ .

*Theorem 5.39.* Let  $(X, \mathcal{T}_X)$  and  $(Y, \mathcal{T}_Y)$  be topological spaces. Then  $(X \times Y, \mathcal{T}_{X \times Y})$  is a topological space.

The Cartesian product of two topological spaces will be assumed to have the product topology unless otherwise specified.

## 5.7 Maps between Topological Spaces - Continuity

We came to the definition of a topology by capturing the notion of “closeness” without using distance. So a topology captures closeness, and hence maps between topological spaces must preserve closeness. In other words, if  $f : X \rightarrow Y$  is a function between topological spaces  $(X, \mathcal{T}_X)$  and  $(Y, \mathcal{T}_Y)$ , then to preserve this structure,  $f$  must send points that are close to points that are close. If  $A$  is a subset of  $X$  and  $p$  a limit point of  $A$ , then  $f(p)$  should be a limit point of  $f(A)$  (or possibly in  $f(A)$ ). We will call such functions continuous.

*Definition.* Let  $(X, \mathcal{T}_X)$  and  $(Y, \mathcal{T}_Y)$  be topological spaces and let  $f : X \rightarrow Y$  be a function. Then  $f$  is *continuous* if for any set  $A \subset X$  and  $p$  a limit point of  $A$ ,  $f(p) \in \overline{f(A)}$ .

You might have guessed that we should define continuous as functions,  $f$ , such that for any open  $U \in \mathcal{T}_X$ ,  $f(U)$  is open, namely in  $\mathcal{T}_Y$ . Though a natural thing to guess, this does not have the desired properties.

*Exercise 5.40.* Let  $f : \mathbb{R}_{std} \rightarrow \mathbb{R}_{std}$  be a constant function. This means that there is a fixed  $y \in \mathbb{R}$  such that  $f(x) = y$  for every  $x$  in  $\mathbb{R}$  (all of  $\mathbb{R}$  is sent to the point  $y$ ). Prove that this function is continuous but that there are open sets  $U$  such that  $f(U)$  is not open. This example shows that requiring  $f(p)$  to be a limit point of  $A$  (rather than in  $\overline{f(A)}$ ) is not a good definition of continuity; explain.

Limit points are hard to work with because of the phrase “for every” in their definition, and the previous definition of continuity is not the one that most mathematicians remember. But it is equivalent to the standard definition, seen in the following theorem.

*Theorem 5.41.* Let  $(X, \mathcal{T}_X)$  and  $(Y, \mathcal{T}_Y)$  be topological spaces and let  $f : X \rightarrow Y$  be a function. Then the following are equivalent.

1. The function  $f$  is continuous,
2. for any set  $V \in \mathcal{T}_Y$ ,  $f^{-1}(V) \in \mathcal{T}_X$ , and
3. for any closed set  $K$  in  $Y$ ,  $f^{-1}(K)$  is closed in  $X$ .

*Exercise 5.42.* Consider the function  $f(x) = mx + b$ , from  $\mathbb{R}_{std}$  to  $\mathbb{R}_{std}$ . Show that  $f$  is continuous for any pair of values  $m$  and  $b$ . You might want to break this exercise into cases depending on whether  $m$  is positive, negative, or zero.

Suppose we have a function  $f : (X, \mathcal{T}_X) \rightarrow (Y, \mathcal{T}_Y)$  from one topological space to another and suppose  $A$  is a subset of  $X$ . It is sometimes convenient to define the restriction map of  $f$  from a subset  $A$  to  $Y$  where we use the relative topology on  $A$  as the topology for  $A$ . So  $f|_A : A \rightarrow Y$  is defined as  $f|_A(a) = f(a)$  for each  $a$  in  $A$ .

*Theorem 5.43.* Let  $f : X \rightarrow Y$  be a function. Suppose  $X = A \cup B$  where  $A$  and  $B$  are closed subsets of  $X$ . If  $f|_A$  is continuous and  $f|_B$  is continuous, then  $f$  is continuous.

As in the group theory chapter, we've discussed what it means to have a map that preserves structure between two topological spaces. Now we're finally ready to say what it means for two topological spaces to be the same.

*Definition* (homeomorphism). A function  $f : X \rightarrow Y$  is a *homeomorphism* if and only if  $f$  is continuous, bijective, and  $f^{-1} : Y \rightarrow X$  is also continuous.

Two topological spaces,  $X$  and  $Y$ , are said to be *homeomorphic* if and only if there exists a homeomorphism  $f : X \rightarrow Y$ .

Some of us are not very comfortable talking about the new function  $f^{-1}$ , the inverse of a bijection  $f$ . The following definition will let us define homeomorphisms only talking about the function  $f$ .

*Definition* (closed and open functions). A continuous function  $f : X \rightarrow Y$  is *closed* if and only if for every closed set  $A$  in  $X$ ,  $f(A)$  is closed in  $Y$ . A continuous function  $f : X \rightarrow Y$  is *open* if and only if for every open set  $U$  in  $X$ ,  $f(U)$  is open in  $Y$ .

*Theorem 5.44.* For a continuous function  $f : X \rightarrow Y$ , the following are equivalent:

- a)  $f$  is a homeomorphism.
- b)  $f$  is bijective and closed.
- c)  $f$  is bijective and open.

*Theorem 5.45.* For numbers  $a < b$  in  $\mathbb{R}^1$ , with the standard topology, the interval  $(a, b)$  is homeomorphic to  $\mathbb{R}^1$ .

*Theorem 5.46.* The letters  $S$  and  $L$ , considered as subsets of  $\mathbb{R}^2$  are homeomorphic.

In group theory, we simply require isomorphisms to be bijective homomorphisms. This forced the inverse to be a homomorphism as well. That does not happen in the topological setting.

*Exercise 5.47.* Prove that there is a continuous bijection from  $\mathbb{R}$  with the standard topology onto  $\mathbb{R}_{LL}$  and yet  $\mathbb{R}$  with the standard topology and  $\mathbb{R}$  with the  $\mathbb{R}_{LL}$  topology are not homeomorphic.

*Theorem 5.48.* The space  $\mathbb{R}^2$  with the standard topology is homeomorphic to the product topology  $\mathbb{R}_{std} \times \mathbb{R}_{std}$ .

The following statement cannot be proven without more rigorous definitions. In what sense could it be made rigorous? Is there a reasonable definition of a “topological property”?

*Metatheorem 5.49.* If  $X$  and  $Y$  are topological spaces and  $f : X \rightarrow Y$  is a homeomorphism, then  $X$  and  $Y$  are the same as topological spaces, *i.e.*, any topological property of the space  $X$  is also a topological property of the space  $Y$ .

## 5.8 Reasons behind the Madness

Although we came to our definition of continuity by requiring functions to send points that are close to points that are close, there are other reasons for the definition that are almost more convincing: all of the simplest functions that we might consider are continuous. Let’s think about this carefully for a few minutes.

*Theorem 5.50.* Let  $(X, \mathcal{T})$  be a topological space. Then  $id_X : X \rightarrow X$ , the identity function, is continuous (actually, it’s a homeomorphism).

*Theorem 5.51.* Let  $(X, \mathcal{T}_X)$  be a topological space and  $A \subset X$  and let  $i_A : A \rightarrow X$  be the inclusion of  $A$  with its subspace topology into  $X$ . Then  $i_A$  is continuous.

Let  $f : Y \rightarrow X$  be a function and  $B \subset Y$ . Then  $f|_B$  is the function formed from  $f$  by restricting the domain of  $f$  to  $B$ .

*Theorem 5.52.* Let  $f : Y \rightarrow X$  be a function. Then  $f$  is continuous if and only if  $f|_{f^{-1}(A)} : f^{-1}(A) \rightarrow X$  is continuous for every set  $A \subset X$ .

*Theorem 5.53.* Let  $f : A \rightarrow X$  be an injection. Then  $f$  is continuous if and only if every set  $B$  of  $A$  for which  $f(B)$  is a relatively open set in  $f(A)$  considered as a subset of  $X$  is open in  $A$ .

*Theorem 5.54.* Let  $X$  and  $Y$  be topological spaces. The projection (first) function  $\pi_X : X \times Y \rightarrow X$  defined by  $\pi_X((x, y)) = x$  is continuous, open, and surjective. Similarly, the (second) projection function  $\pi_Y : X \times Y \rightarrow Y$  defined by  $\pi_Y((x, y)) = y$  is continuous, open, and surjective. For any topology on the Cartesian product that omits any open set in the product topology, at least one of the projection maps will not be open.

*Exercise 5.55.* Let  $X$  and  $Y$  be topological spaces. The projection function  $\pi_X : X \times Y \rightarrow X$  need not be a closed map.

*Theorem 5.56.* Let  $X$ ,  $Y$ , and  $Z$  be topological spaces and suppose we have continuous functions  $f : Z \rightarrow X$  and  $g : Z \rightarrow Y$ . Then there is a unique continuous function  $h : Z \rightarrow X \times Y$  such that  $\pi_X \circ h = f$  and  $\pi_Y \circ h = g$ . Furthermore, a function  $k : Z \rightarrow X \times Y$  is continuous if and only if  $\pi_X \circ k$  and  $\pi_Y \circ k$  are both continuous.

Actually, this whole thing makes more sense backwards. The definition of continuity is what it has to be to preserve closeness, and then the subspace and product topologies are defined so that the related functions (the inclusions and projections) described above are continuous.

## 5.9 Connected Spaces

A topology somewhat captures the intuitive idea of closeness, which is a local property; however the topology can give us some indication of global properties of the space as well. Think of one of those automatic cleaner robots that keeps going unless it hits a wall. If you turn it on and leave it on the ground floor of your house, it will never clean upstairs. The set that is the floor space in your house has at least two distinct pieces. Thinking of this set with a topology like the standard topology, we would say that it is disconnected. So one global feature of a space that is captured by a topology is a sense of connectivity.

Following our intuition about the two story house, it turns out to be easier to say what we mean by a space not being connected and then we define being connected by saying it isn't not connected.

*Definition.* A topological space  $(X, \mathcal{T})$  is called *connected* if it cannot be written as the disjoint union of two *non-empty* open sets,  $U, V \in \mathcal{T}$  with  $U \cap V = \emptyset$  and  $X = U \cup V$ . In other words, a space is not connected if it's made of two non-touching chunks that are open sets.

For example, consider the set  $Y = [-2, -1) \cup \{1, 3, 5\} \subset \mathbb{R}$ . If we use the standard topology on  $\mathbb{R}$  and consider the subspace topology on  $Y$ , then  $Y$  is disconnected. We can see this fact as follows. Recall that  $U = (-3, -1) \in \mathcal{T}_{std}$  and  $V = (0, 6) \in \mathcal{T}_{std}$ . But then  $W = U \cap Y = [-2, -1)$  and  $Z = V \cap Y = \{1, 3, 5\}$  are open in  $\mathcal{T}_Y$ . But  $Y = W \cup Z$ ,  $W \cap Z = \emptyset$ , and clearly neither  $W$  nor  $Z$  is empty.

*Exercise 5.57.* Consider  $\mathbb{R}$  with our six favorite topologies: standard, discrete, indiscrete, finite complement, countable complement, and lower limit. For which topologies is  $\mathbb{R}$  a connected topological space and for which is it not? Explain your answers.



*Theorem 5.58.* Let  $(X, \mathcal{T}_X)$  and  $(Y, \mathcal{T}_Y)$  be topological spaces and let  $f : X \rightarrow Y$  be a surjective continuous function. If  $X$  is connected, then  $Y$  is connected.

## 5.10 Metric Topologies and Continuity

We began this whole inquiry by attempting to remove distance from the statement of some geometric questions. But having a sense of distance actually led us to the standard topology. Does this happen in general? Yes. If  $X$  is a set with a distance, there is a topology on  $X$  such that, for any subset  $A$ , the limit points of  $A$  in the topology are arbitrarily close to  $A$  using the distance.

First we must be careful about what we mean by a sense of distance. A sense of distance is actually a function that takes as input two points in the set and produces a positive real number, which we consider the distance between the two points. Notice that for distance to be a coherent concept, it must be shorter to go from point A to B than it is to go from A to C and then to B. This relationship among distances is called the triangle inequality.

*Definition.* A *metric* on  $X$ , is a function  $d : X \times X \rightarrow \mathbb{R}$  such that

1.  $d(x, y) \geq 0$  for all  $x, y \in X$ ,
2.  $d(x, y) = 0$  if and only if  $x = y$ , and
3. for any three points  $x, y, z \in X$ ,  $d(x, z) \leq d(x, y) + d(y, z)$ .

For example, the distance between points in Texas “as the crow flies” is a metric. Certainly the distance between any two points is non-negative. The distance between two places is 0 if and only if they are the same place. And it is certainly shorter for a crow to fly directly from Dallas to Houston than it is to fly Dallas to Austin and then Austin to Houston.

*Theorem 5.59.* Consider the function  $d : \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$  defined by

$$d((x_1, x_2), (y_1, y_2)) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}.$$

Then  $d$  is a metric, called the Euclidean metric.

This metric actually generalizes to  $\mathbb{R}^n$  for any  $n$ . If  $x = (x_1, \dots, x_n)$  and  $y = (y_1, \dots, y_n)$  are two points in  $\mathbb{R}^n$ , then  $d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$  is a metric on  $\mathbb{R}^n$ , also called the Euclidean metric. Notice that on  $\mathbb{R}^1$ ,  $d(x_1, y_1) = \sqrt{(x_1 - y_1)^2} = |x_1 - y_1|$ .

Once we have a metric, we have the collection of points that are within a certain distance of a point: discs and spheres, etcetera, which we will call balls.

*Definition.* Let  $p \in \mathbb{R}^n$  and  $\varepsilon > 0$  a real number and  $d$  the Euclidean metric. Then we define the ball centered at  $p$  of radius  $\varepsilon$  centered at  $p$  as

$$B(p, \varepsilon) = \{q \in \mathbb{R}^n \mid d(p, q) < \varepsilon\}.$$

And now for the topology induced by the metric.

*Example 7* (metric topology). Let  $X$  be a set and  $d$  a metric on  $X$ . Then let  $\mathcal{T}_d$  be defined as follows. A subset  $U \subset X$  is in  $\mathcal{T}$  if and only if for every point  $p \in U$ , there is an  $\varepsilon_p > 0$  such that  $B(p, \varepsilon_p) \subset U$ . In other words, these balls are a basis for the topology  $\mathcal{T}_d$ , which is called the *metric topology* on  $X$ .

*Theorem 5.60.* The Euclidean metric topology on  $\mathbb{R}^n$  is a topology. Furthermore,  $\mathbb{R}^n$  with the metric topology is homeomorphic to  $\mathbb{R}^n$  with the standard topology.

In the sections above, we are claiming to measure whether a point  $p$  is close to a set  $A$  using limit points. We need to be able to talk about the distance between points and sets to compare our theorems.

*Definition.* Let  $X$  be a set and let  $d$  be a metric on  $X$ . Also, let  $p$  be a point in  $X$  and  $A \subset X$ . Then the distance between  $p$  and  $A$ , which we will write, abusing notation, as

$$d(A, p) = \inf \left( \{d(a, p) \mid a \in A\} \right).$$

The symbol  $\inf$  is pronounced “infimum” and is the biggest lower bound for a collection of real numbers. For example,

$$\inf \left( [1, 6) \right) = 1, \quad \inf \left( \left( \frac{1}{2}, 3 \right] \right) = \frac{1}{2}, \quad \text{and} \quad \inf \left( \left\{ \frac{1}{n} \mid n \in \mathbb{N} \right\} \right) = 0.$$

*Theorem 5.61.* Let  $X$  be a set with a metric,  $d$ . Also, let  $p \in X$  and  $A \subset X$ . Then  $d(A, p) = 0$  if and only if  $p \in \overline{f(A)}$  in the metric topology on  $X$ .

Sometimes the same topology can be described using a metric and alternatively described in a different way.

*Definition.* Let  $X$  be a set and consider the function  $d_0 : X \times X \rightarrow \mathbb{R}$  defined by  $d_0(x, y) = 0$  if  $x = y$  and  $d_0(x, y) = 1$  if  $x \neq y$ . We call  $d_0$  the trivial metric.

*Exercise 5.62.* Show that  $d_0$  defined above is a metric. The metric topology induced by  $d_0$  is already familiar to us; what is its name?

Let us now review the definition of continuity that you probably first encountered in calculus—the  $\varepsilon$ - $\delta$  definition.

*Definition* (analytic continuity in  $\mathbb{R}$ ). A function  $f : D \subseteq \mathbb{R} \rightarrow \mathbb{R}$  is *continuous at  $x$*  (with respect to the metric  $d$ ) if and only if for every  $\varepsilon > 0$  there is a  $\delta > 0$  such that if  $z \in D$  and  $|x - z| = d(x, z) < \delta$  then  $|f(x) - f(z)| = d(f(x), f(z)) < \varepsilon$ . We say  $f$  is (analytically) *continuous* if it is continuous for every point  $x$  in its domain  $D$ .

And just so that it's in the same section, we'll recall the equivalent definitions of continuity in a topology.

*Definition* (topological continuity). A function  $f : (X, \mathcal{T}_X) \rightarrow (Y, \mathcal{T}_Y)$  is *continuous* if for every set  $V \in \mathcal{T}_Y$ ,  $f^{-1}(V) \in \mathcal{T}_X$ .

*Theorem 5.63.* A function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is topologically continuous in the standard topology if and only if it is analytically continuous with respect to the Euclidean metric.

The previous theorem is actually a special case of the following theorem.

*Theorem 5.64.* If  $X$  and  $Y$  are metric spaces with metrics  $d_X$  and  $d_Y$  respectively, then a function  $f : X \rightarrow Y$  is analytically continuous with respect to these metrics if and only if it is topologically continuous with respect to the metric topologies.

## 5.11 Abstraction is Useful

“The Unreasonable Effectiveness of Mathematics in the Natural Sciences” is the name of a famous paper by Eugene Wigner, which makes the case suggested in the title. In the case of topology, we generated the definitions and the ideas by abstracting concepts from other mathematics that was already an abstraction of reality. We might well imagine that we had strayed rather far from the real world by this point. However, the perspectives that emerge from a study of topology actually have direct applicability to questions about our world. Modern theories about the shape of the universe often are described using topological spaces. Some features of the structure of DNA molecules are illuminated through the study of topology. Analysis of computing networks involves topology.

The step of abstracting ideas to build coherent mathematical structures has been one of the most fruitful lines of intellectual exploration in history.

Graph theory, group theory, calculus, and topology all illustrate the rich world of ideas that we can explore when we liberate ourselves to develop mathematical ideas in the abstract. But all these areas have also demonstrated their uncanny capacity for shedding light on the secrets of our real world as well.

## Appendix A

# Appendix: Sets and Functions

Essentially all modern mathematical structures involve sets and functions. Certainly graphs, groups, analysis, and topology all discuss sets with various properties and functions that relate these sets. This appendix presents the basic vocabulary used to talk about sets and functions.

### A.1 Sets

A **set** is a collection of **elements** (think “things” or “points”). There are two standard ways of describing a set: listing the elements or describing a property that defines the elements. The key property is that an element is either definitely in or definitely not in a given set.

Example (1): Zero is in the set  $X$  that contains only zero, two, three, and five, written as

$$0 \in X = \{0, 2, 3, 5\}.$$

The symbol “ $\in$ ” is read as “is in” or “is an element of”. Also, we use a pair “{” and “}” around the list or description of the elements of a set.

Example (2): Similarly, zero is not in the set  $Y$  of odd integers, written

$$0 \notin Y = \{m \mid m \text{ is an odd integer}\}.$$

The symbol “ $\notin$ ” is read “is not in” or “is not an element of”.

You may have noticed the symbol “ $\mid$ ” in the description of  $Y$ . This symbol is read “such that”, and it means that  $Y$  is the collection of objects before the “ $\mid$ ” subject to the conditions described after it.

Example (3): So, we can describe the even integers as

$$\{n | n \text{ is an integer and there exists an integer } m \text{ such that } n = 2m\}.$$

There are many sets that appear so often that we want to have a compact way of writing them. So here are some commonly appearing sets and the symbols we will use to reference them:

*Notation.* 1. The set  $\{1, 2, 3, 4, 5, \dots\}$  is called the **natural numbers** and is written  $\mathbb{N}$ .

2. The set  $\{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\}$  is called the **integers** and is written  $\mathbb{Z}$ .

3. The set of all decimal numbers, both positive and negative, is called the **real numbers** and written  $\mathbb{R}$ .

4. The set of all ratios of integers  $\frac{a}{b}$ , where the denominator is non-zero and two ratios are the same if they are equal when put in lowest terms is called the **rational numbers**, written

$$\mathbb{Q} = \left\{ \frac{a}{b} \mid a, b \in \mathbb{Z}, b \neq 0, \frac{a}{b} = \frac{c}{d} \text{ if and only if } ad - bc = 0 \right\}.$$

5. The set containing no elements is called the **empty set** and written  $\emptyset$ .

When one set  $A$  is contained inside a second set  $B$ , we say that  $A$  is a **subset** of  $B$ , written  $A \subset B$ . If  $A$  is a subset of  $B$  that does not contain all of  $B$ , then we may write  $A \subsetneq B$  to denote that  $A$  is a proper subset of  $B$ . For example,  $\mathbb{N} \subset \mathbb{Z}$  and  $\emptyset \subsetneq \mathbb{Q} \subsetneq \mathbb{R}$ .

If  $A$  is a subset of  $B$ , we may want to describe the elements that are in  $B$  but *not* in  $A$ . We write this set as

$$B \setminus A = \{b \in B \mid b \notin A\}.$$

This slash mark is called *setminus*, and it means to remove the second set from the first. Think of it as the word “except”. For example, if we want to talk about the non-zero real numbers (all of  $\mathbb{R}$  except  $\{0\}$ ), we will use the symbol  $\mathbb{R} \setminus \{0\}$ . Similarly, the non-zero rational numbers will be denoted by  $\mathbb{Q} \setminus \{0\}$ . Also, the *irrational numbers* are the real numbers that are not rational, written as  $\mathbb{R} \setminus \mathbb{Q}$ .

Note the difference between being a subset and being an element of a set. Zero is an element of the integers, but the set containing only zero is a subset of the integers.

$$0 \in \mathbb{Z} \quad \text{vs.} \quad \{0\} \subset \mathbb{Z}$$

Given two sets,  $A$  and  $B$ , we can build several new sets from them.

*Notation.* Let  $A$  and  $B$  be two sets.

1. The set of elements contained in at least one of  $A$  or  $B$  is called the **union** of  $A$  and  $B$  and written

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

2. The set of elements that are in both  $A$  and  $B$  is called the **intersection** of  $A$  and  $B$  and written

$$A \cap B = \{x | x \in A \text{ and } x \in B\}.$$

3. The set of ordered pairs with first element in  $A$  and second element in  $B$  is called the **Cartesian product** of  $A$  and  $B$  and written

$$A \times B = \{(a, b) | a \in A, b \in B\}.$$

For example,

$$\{1, 2, 3, 4, 6\} \cup \{2, 3, 5, 6, 8\} = \{1, 2, 3, 4, 5, 6, 8\},$$

$$\{1, 2, 3, 4, 6\} \cap \{2, 3, 5, 6, 8\} = \{2, 3, 6\}, \text{ and}$$

$$\{1, 2\} \times \{2, 5\} = \{(1, 2), (1, 5), (2, 2), (2, 5)\}.$$

If  $A' \subset A$  and  $B' \subset B$ , then we can think of  $A' \times B'$  as a subset of  $A \times B$ . So  $\{1, 2\} \times \{2, 5\} \subset \mathbb{N} \times \mathbb{N}$ .

Instead of  $A \times A$  we will sometimes write  $A^2$ . For example, you are probably familiar with the collection of ordered pairs of real numbers,  $\mathbb{R} \times \mathbb{R} = \mathbb{R}^2$ .

We will often use lowercase letters for elements and uppercase letters for sets, though not always. It is a good idea to have related symbols for a set and an element of the set. For example,  $n$  will usually denote a natural number, in which case  $n \in \mathbb{N}$ . Similarly, if  $G$  is a set,  $g$  will usually denote an element of  $G$ . Whenever choosing notation for a theorem or proof, try to make it clear how the symbols relate, possibly by using this convention.

## A.2 Functions

We've pinned down the notion of a set; now all we need is a good definition of a function. Modern mathematics uses functions to study almost everything, and the word “function” is even used in common speech, so we must take extra care.

You are familiar with functions such as  $f(x) = x^2$  and  $g(t) = 3\sqrt{t} - 17$  that take a real number  $x$  or  $t$ , do something to it, and produce another real number. Sometimes a function is described as a rule, formula, or equation that takes a number (the independent variable) and gives another (the dependent variable). However, this description is a little vague and a little restrictive; it is neither specific enough to be precise nor general enough to satisfy our needs for developing more abstract mathematical ideas. So, once again, we will use the intuition that we have developed from our common experience with functions and then abstract those ideas to create a more formal definition of function.

Familiar functions such as  $f(x) = x^2$  take one number and produce another number. But we can generalize that idea to view a function as relating each element of one set with an element of another set. And instead of some rule being used to make the relation, any specified relation will be called a function.

A function  $f$  associates each element of a set  $D_f$ , the **domain** of  $f$ , with an element of a set  $C_f$ , the **codomain** of  $f$ . For example, if the domain is  $D_f = \{\clubsuit, \diamond, \heartsuit, \spadesuit\}$  and the codomain is  $C_f = \{a, b, c\}$ , then associating  $\clubsuit \mapsto b$ ,  $\diamond \mapsto a$ ,  $\heartsuit \mapsto b$ , and  $\spadesuit \mapsto a$  is a function. Since a function relates each element of  $D_f$  with an element in  $C_f$ , we can view a function as a collection of ordered pairs, as follows.

*Definition.* A **function**,  $f$ , from domain set  $D_f$  to codomain set  $C_f$  is a subset of the Cartesian product  $D_f \times C_f$  such that each element of  $D_f$  occurs in exactly one of those pairs. Sometimes a function  $f$  is written as  $f : D_f \rightarrow C_f$ , where  $f(d) = c$  means that  $(d, c)$  is an element of the function  $f \subset D_f \times C_f$ .

There are several different notations for functions that we will use, so we should compare them. Let  $f \subset D_f \times C_f$  be a function.

(A): Thinking about  $f$  as a collection of ordered pairs is like considering the graph of the kinds of functions you studied in highschool. The requirement that each element of  $D_f$  appear in exactly one ordered pair is a reformulation of the “Vertical Line Test”. Of course, you can only draw the graph of a function when the domain and codomain are subsets of the



reals. So, sometimes we will think of functions like “graphs”.

(B): Sometimes we will think of  $f$  as a map from  $D_f$  to  $C_f$ . For example, if  $f$  contains the ordered pair  $(d, c)$ , then we think of it as sending  $d$  to  $c$ , which we could write as  $d \mapsto c$  or  $d \xrightarrow{f} c$ . So sometimes we will think of functions as maps that send domain elements to codomain elements.

(C): We could also think of  $f$  as a rule that associates an element of  $C_f$  with each element of  $D_f$ . When thinking of our function like that, we will instead write  $(d, c) \in D_f \times C_f$  as  $f(d) = c$ . If  $f$  is described in terms of a formula, then we will think of it as a rule that relates domain elements to codomains elements.

These notations are all conveying the same information:

$$(d, c) \in D_f \times C_f \text{ if and only if } d \xrightarrow{f} c \text{ if and only if } f(d) = c.$$

If the domain and codomain are familiar sets, like subsets of  $\mathbb{R}$ , then sometimes the function is given by a formula. For example, the relation that sends a “number” to “three times the number squared” is a function from  $\mathbb{R}$  to  $\mathbb{R}$ . If we call this function  $f$ , we often write  $f : \mathbb{R} \rightarrow \mathbb{R}$  is defined by  $f(x) = 3x^2$ . If we viewed this function as a subset of the Cartesian product  $\mathbb{R} \times \mathbb{R} = \mathbb{R}^2$  it would look like a parabola, and we would write  $f = \{(x, 3x^2) | x \in \mathbb{R}\}$ .

*Exercise A.1.* For each of the following relations, determine if it is a function.

1.  $\{(\text{person}, \text{his/her birthday})\} \subset \text{People} \times \text{Dates}$
2.  $\{(\text{person's first name}, \text{his/her last name})\} \subset \text{Names} \times \text{Names}$
3.  $\{(\text{natural number}, \text{its remainder when divided by 3})\} \subset \mathbb{N} \times \{0, 1, 2\}$

A single function, by itself, is not that interesting; we need to know how to combine functions. If  $f : A \rightarrow B$  and  $g : B \rightarrow C$  are functions, then we want to be able to combine them into a new function  $h : A \rightarrow C$  by first doing  $f$  then doing  $g$ . For example, if  $f$  is the function “ $x \mapsto 4x^2$ ” from  $\mathbb{R}$  to  $\mathbb{R}$  and  $g$  is the function “ $t \mapsto \frac{t}{2} + 1$ ” from  $\mathbb{R}$  to  $\mathbb{R}$ , then combining them in this way produces a function that maps  $x$  to  $\frac{(4x^2)}{2} + 1 = 2x^2 + 1$ . This process is called *composition* and is written  $h = g \circ f : A \rightarrow C$ . Note the order of the component functions,  $f$  and  $g$ . The composition  $h = g \circ f$  means to perform  $f$  first and then perform  $g$  to the result. The reason for this choice of order in the notation is that  $h(x) = g(f(x))$ . The  $f$  is closer to the domain element because it is the first function performed.

In the above example,  $C_f = D_g$ , which guarantees that  $h = g \circ f$  makes sense. However, this condition is not a necessary condition for the composition to exist. For their composition to make sense, all we really need is that all outputs from  $f$  are possible inputs for  $g$ .

*Definitions.* Let  $f : D_f \rightarrow C_f$  be a function.

1. Define the **range** of  $f$  as

$$R_f = \{c \in C_f \mid \text{there exists a } d \in D_f \text{ such that } f(d) = c\}.$$

2. Similarly, if  $D' \subset D_f$ , then  $f(D') = \{c \in C_f \mid \text{there exists a } d \in D' \text{ such that } f(d) = c\}$  is called the **image** of  $D'$ . So  $R_f$  is the image of the entire domain,  $f(D_f)$ .
3. And if  $C' \subset C_f$ , then  $f^{-1}(C') = \{d \in D_f \mid f(d) \in C'\}$  is called the **preimage** of  $C'$ .

*Definition.* Let  $f : D_f \rightarrow C_f$  and  $g : D_g \rightarrow C_g$  be functions such that  $R_f \subset C_g$ . Then define a new function, called the **composition** of  $f$  and  $g$ ,  $g \circ f : D_f \rightarrow C_g$  by  $x \xrightarrow{g \circ f} g(f(x))$ .

### A.3 Special Functions

There are a number of special kinds of functions that you might see when reading mathematics. Here are a few.

*Notation.* Let  $X$  and  $Y$  be sets.

1. There is always a function from  $X$  to  $X$ . Define the **identity function** on  $X$ ,  $Id_X : X \rightarrow X$  by  $x \mapsto x$  for each  $x \in X$ .
2. If  $Y \subset X$ , then we can define the **inclusion** of  $Y$  into  $X$ ,  $i_Y : Y \rightarrow X$  by  $y \mapsto y$ . So the identity function is the inclusion of a set into itself.
3. There are two functions naturally associated with  $X \times Y$ . Let  $\pi_X : X \times Y$  be defined by  $\pi_X((x, y)) = x$  and  $\pi_Y : X \times Y \rightarrow Y$  by  $\pi_Y((x, y)) = y$ . Then  $\pi_X$  and  $\pi_Y$  are called the **projections** onto the first and second factors respectively.
4. Suppose  $f : Y \rightarrow X$  is a function and  $Z \subset Y$ . Then we can think of  $f$  as a function  $f|_Z : Z \rightarrow X$  by defining  $f|_Z(z) = f(z)$ . We will call  $f|_Z$  the **restriction** of  $f$  to  $Z$ . In other words,  $f \subset Y \times X$ , and  $Z \subset Y$  so that  $(Z \times X) \subset (Y \times X)$ , then  $f|_Z = f \cap (Z \times X)$ . Sometimes we will abuse notation and write  $f|_Z$  as  $f$ .

A function relates two sets, its domain and codomain. In general, not every element of the codomain is associated to an element in the domain. The collection of elements in the codomain of a function  $f$  that are associated with domain elements is called the **range** or **image** of  $f$ , denoted  $R_f$  or  $Im(f)$ . The range is a subset of the codomain. A function is called **surjective** or **onto** if every element of the codomain is associated to at least one element of the domain. In other words, a function is surjective if and only if its range equals its codomain. A function is called **injective** or **1-to-1** if no element in the range/codomain is associated to more than one element in the domain. A function that is both injective and surjective is called **bijective**.

*Exercise A.2.* Let  $f(x) = \frac{3}{x^2}$ , which is a function from  $\mathbb{R} \setminus \{0\}$  to  $\mathbb{R}$ . What is the range of  $f$ ,  $R_f$ ?

*Exercise A.3.* A polynomial  $p(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$  is a function from  $\mathbb{R}$  to  $\mathbb{R}$ , given by plugging the domain number into the polynomial as  $x$ .

1. Give an example of an injective polynomial.
2. Give an example of a polynomial that is not injective.
3. Give an example of a surjective polynomial.
4. Give an example of a polynomial that is not surjective.
5. Give an example of a bijective polynomial.
6. Give an example of a non-bijective polynomial.
7. (optional) Suppose  $f(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$ . Find conditions on the coefficients that determine whether  $f$  is injective, whether  $f$  is surjective, and whether  $f$  is bijective.

Some properties of injectivity and surjectivity are preserved by composition of functions, and some are not.

*Theorem A.4.* Let  $f : D_f \rightarrow C_f$  and  $g : D_g \rightarrow C_g$  be functions such that  $R_f \subset D_g$ .

1. If  $f$  and  $g$  are injective, then  $g \circ f$  is injective.
2. If  $f$  and  $g$  are surjective, then  $g \circ f$  is surjective.

3. If  $f$  is injective, then there is a function,  $h : R_f \rightarrow A$  such that  $h \circ f : A \rightarrow A$  is the identity function on  $A$ .

*Exercise A.5.* Give an example of a pair of functions  $f : A \rightarrow B$  and  $g : B \rightarrow C$  such that  $g$  is not injective but  $g \circ f$  is injective.

If  $B'$  is a subset of the codomain of a function,  $f : A \rightarrow B$ , that contains the range, then we can think of  $f$  as a function  $f : A \rightarrow B'$  as well. Altering a function in this way may change many things about it. For example, if  $B' = R_f$ , then the new function is surjective. When we change the domain or codomain of a function, we must argue that the change does not affect the answer to the problem we are considering.

## A.4 Binary Operations

Note that the domain of a function does not have to be a simple set; it could even be a Cartesian product itself. For example, the rule  $f(x, y, z) = 3xy - 17z + y^2$  is a function from  $\mathbb{R}^3 = \mathbb{R} \times \mathbb{R} \times \mathbb{R}$  to  $\mathbb{R}$ . For a more familiar example,  $g(x, y) = x + y$  is a function from  $\mathbb{R} \times \mathbb{R}$  to  $\mathbb{R}$ . This function is exactly one of our generative binary operations in the Chapter *Group Theory*. And so the concepts of Cartesian product and function let us define precisely what we mean by a binary operation on a set.

*Definition.* If  $S$  is a set, then a **binary operation** on  $S$  is a function  $* : S \times S \rightarrow S$ . This function is actually a subset of  $(S \times S) \times S$  such that every pair  $(s_1, s_2)$  appears in exactly one element  $((s_1, s_2), s_3)$  in the subset. A set  $S$  together with a binary operation  $*$  on  $S$  is written  $(S, *)$ .

A few examples of sets with binary operations that you're familiar with are  $(\mathbb{Z}, +)$  and  $(\mathbb{R}, \cdot)$ , the integers with addition and the real numbers with multiplication, respectively. There are many others. And just as we write,  $a + b$  instead of  $+((a, b))$ , we cheat and write  $s_1 * s_2$  to denote the value of the  $*$  function applied to the element of the domain  $(s_1, s_2)$ .

If  $s_1 * s_2$  is not an element of  $S$  for some pair  $(s_1, s_2) \in S \times S$  then  $*$  cannot be a binary operation, because it is not a function. When  $s_1 * s_2$  is not an element of  $S$ , we say that  $*$  is not *closed*. And if  $*$  is given by a rule (instead of a set of ordered pairs) then sometimes  $s_1 * s_2$  gives two different values, depending how we compute it, then  $*$  is not a binary operation because again it is not a function. When  $s_1 * s_2$  gives two different values, depending how we compute it; we say that  $*$  is not *well-defined*. By requiring binary operations to be functions, we require them to be closed and well-defined.

A binary operation  $*$  on a set  $S$  is called **associative** if for every  $a, b$ , and  $c$  in  $S$ ,  $(a * b) * c = a * (b * c)$ . Parentheses mean what they always have: do the operation on the inside first. Associativity is often quite hard or tedious to verify, but it guarantees that we don't have to write parentheses when we perform several operations in a row, which is well worth the effort.

*Fact 1.* Addition and multiplication of real-valued functions (functions whose codomains are subsets of  $\mathbb{R}$ ) are all associative binary operations. Composition of functions is associative, as long as it makes sense to compose them.

*Exercise A.6.* Show that the following binary operations are or are not associative.

1.  $a * b = a$  on any set
2.  $a * b = a - b$  on  $\mathbb{Q}$

The restriction of a binary operation is not necessarily a binary operation. If  $*$  is a binary operation on  $S$ , then it is a function from  $S \times S$  to  $S$ . Restricting it to a subset  $T \subset S$  produces a function  $T \times T \rightarrow S$ , not  $T \times T \rightarrow T$  as we would need. We give a definition for the situation in which the restriction actually is a binary operation.

*Definition.* If  $(S, *)$  is a set with a binary operation and  $T$  is a subset of  $S$ , then we say that  $T$  is **closed under  $*$**  if, for all pairs in  $T \times T$ ,  $t_1 * t_2$  is an element of  $T$  (as opposed to just in  $S$ ). In other words, the restriction  $* : T \times T \rightarrow S$  has range inside  $T$  and thus may be considered as a binary operation on  $T$ .

We should write this restriction as  $*|_{T \times T}$ , but we will usually just write  $*|_T$ . Note that, if  $*|_T$  is a binary operation that is the restriction of an associative binary operation  $*$ , then  $*|_T$  is associative on  $T$ . We say that subsets *inherit* associativity from the set.

*Exercise A.7.* Consider the non-zero reals with the binary operation of multiplication,  $(\mathbb{R} \setminus \{0\}, \cdot)$ . Are the non-zero rational numbers,  $\mathbb{Q} \setminus \{0\}$ , closed under multiplication? Are the irrational numbers,  $\mathbb{R} \setminus \mathbb{Q}$ , closed under multiplication? Justify your answer.

## A.5 Cardinality

Modern mathematics uses functions to answer many questions that are hard to think about directly. For example, it is clear when two finite sets have the same “size”; they just have the same number of elements. But what

does it mean for two infinite sets  $A$  and  $B$  to have the “same size”? We may not know the answer, but certainly if  $A$  sits inside  $B$ , then  $B$  should be at least big as  $A$ .

*Definition.* Let  $A$  be a set. Define the **cardinality** of  $A$ , written  $|A|$  as the number of elements in  $A$ . If  $A$  has a finite number of elements, then  $|A|$  is that finite number. If  $A$  has an infinite number of elements, then  $|A| = \infty$ .

The previous definition makes it seem like all infinite sets are the same size. But this is not true.

*Definition.* If  $A$  and  $B$  are sets, then we say that  $|A| \leq |B|$  if there is an injection  $f : A \rightarrow B$ . We say  $|A| = |B|$  if there is a bijection from  $A$  to  $B$ .

In particular,  $|\mathbb{N}| \leq |\mathbb{Z}| \leq |\mathbb{Q}| \leq |\mathbb{R}|$  because the inclusion maps are injective. All four of these are infinite sets, but are they all the same size?

*Theorem A.8 (Cantor).* There does not exist a surjection  $f : \mathbb{N} \rightarrow \mathbb{R}$ . However, there does exist a bijection  $g : \mathbb{N} \rightarrow \mathbb{Q}$ .

*Definition.* Let  $A$  be a set. If there is a surjection from (some subset of)  $\mathbb{N}$  onto  $A$ , then we say that  $A$  is **countable**. If not, then we say that  $A$  is **uncountable**.

So Cantor’s Theorem tells us that  $|\mathbb{R}|$  is strictly bigger than  $|\mathbb{N}|$ . This fact doesn’t even make sense without defining cardinality using functions.

## A.6 Notation

$\in$   
 $\{a, b\}$   
 $\setminus$   
 $\cup, \cap$   
 $\mathbb{N}, \mathbb{Z}, \mathbb{R}, \mathbb{Q}$   
 $Id_X$   
 $|A|$

# Index

- $\mathbb{R}_{LL}$ , 124
- basic open set, 123
- basis element, 123
- basis of a topology, 123
- closed function, 126
- closed set, 119
- closure, 119
- co-finite topology, 117
- discrete topology, 116
- finite complement topology, 117
- function
  - closed, 126
  - open, 126
- homeomorphism, 126
- indiscrete topology, 115
- integer, 134
- isolated point, 118
- limit point, 114
- lower limit topology, 124
- metric topology, 130
- natural number, 134
- open function, 126
- open set, 114
- Product, 124
- product topology, 124
- set
  - basic, 123
  - closure of, 119
  - open, 114
- Sorgenfrey line, 124
- standard topology on  $\mathbb{R}^2$ , 116
- standard topology on  $\mathbb{R}^n$ , 117
- subspace topology, 122
- topological space, 114
  - homeomorphic, 126
- topologist's comb, 121
- topologist's sine curve, 121
- topology
  - basis, 123
  - co-finite, 117
  - discrete, 116
  - finite complement, 117
  - indiscrete, 115
  - lower limit, 124
  - standard
    - $\mathbb{R}^2$ , 116
    - $\mathbb{R}^n$ , 117
  - subspace, 122

Brian Katz  
Department of Mathematics  
RLM 8.100  
The University of Texas at Austin  
Austin, TX 78712  
bkatz@math.utexas.edu

Michael Starbird  
Department of Mathematics  
RLM 8.100  
The University of Texas at Austin  
Austin, TX 78712  
starbird@math.utexas.edu