

MATH 240 Fall 2024

notes by
MARCEL GOH

A note on these notes. After each class, this document will be updated with the new material that was just covered. The datestamps in the left margin indicate when the notes from each day start. Subsections labelled with a * are optional. This document is heavily based on notes by Jeremy Macdonald, but any errors are likely my own. Please email me if you find any.

I. Foundations	2
1. Set theory	3
2. Propositional logic	10
3. Predicate logic	17
4. Proofs	21
5. Functions	28
6. Cardinality	33
7. Relations	37
II. Number theory	41
8. Division	42
9. Primes	47
10. Modular arithmetic	51
11. Applications of number theory	60
III. Graph theory	68
12. Definitions and basic notions	69
13. Triangles and bipartite graphs	74
14. Trees	78
15. Eulerian trails and circuits	81
16. Planar graphs	83
IV. Combinatorics	88
17. Basic counting techniques	89
18. Permutations and combinations	92
19. Applications of the pigeonhole principle	100
20. Recurrences	102

I. FOUNDATIONS

*Die Mathematik ist in ihrer Entwicklung völlig frei
und nur an die selbstredende Rücksicht gebunden,
dass ihre Begriffe sowohl in sich widerspruchlos sind,
als auch in festen durch Definitionen geordneten Beziehungen
zu den vorher gebildeten,
bereits vorhandenen und bewährten Begriffen stehen.*

— GEORG CANTOR, *Grundlagen einer allgemeinen Mannigfaltigkeitslehre* (1883)

1. Set theory

29.VIII A *set* is a collection of distinct objects, called its *elements* or its *members*. If x is a member of set A , then we write $x \in A$, and if x is not an element of A , then we write $x \notin A$. Sets can be written by listing out its elements. For example,

$$\{1, 4, 7, 10, \sqrt{782}\}$$

and

$$\{\{1, 2\}, \pi, \{4\}\}$$

are both sets (the second example shows that sets can themselves contain other sets). The order of the elements is not important, and duplicate elements are ignored (so $\{1, 2\} = \{2, 1\} = \{1, 1, 2\}$). The notation using $\{$ and $\}$ is useful for defining small, concrete examples, but expressing large sets can become very cumbersome. The first way one can describe larger sets is to use the \dots symbol and the power of suggestion. For instance, anyone faced with the notation

$$A = \{1, 3, 5, 7, 9, \dots\}$$

can quickly guess that this set is supposed to contain all the positive odd integers. We can also use \dots to define finite sets. Most Canadians will be able to tell you that the set

$$\{\text{Alberta, British Columbia}, \dots, \text{Yukon}\}$$

of provinces and territories contains 13 elements. But this notation inherently produces some ambiguity. For example, since the sequence of positive palindromic binary numbers starts 1, 3, 5, 7, 9, 15, 17, 21, 27, \dots , we are left with some doubt as to whether the set A above should be the set of odd numbers or the set of palindromic binary numbers.

But there is another, less ambiguous way to define large sets. It is called set-builder notation and it refers to any construction of the form

$$\{x \in U : P(x)\},$$

where x is a variable, U is a set, and P is a statement about x . The resulting set contains *all x such that $P(x)$ holds*. For example, letting \mathbf{N} denote the set $\{0, 1, 2, 3, \dots\}$ of counting numbers (more on this later), to define the set of all odd numbers, we can write

$$A = \{x \in \mathbf{N} : \text{there exists } k \in \mathbf{N} \text{ such that } x = 2k + 1\}.$$

Note that the statement $P(x)$ must contain x , but it may also contain other previously defined symbols, as well as new symbols defined within the statement (such as k in the example above).

Special sets of numbers. There are certain infinite sets of numbers that are used so often as to be given special bold notation. Back in elementary school,

you learned to use the counting numbers $0, 1, 2, 3, \dots$. We already saw this set in the previous paragraph; in the business, this set is known as the *natural numbers*, because if you go on a nature hike you can use them to count the number of bluebells, donkeys, etc. that you see. (Many mathematicians use the symbol \mathbf{N} to denote this set without zero. When in doubt, clarify with the person you're talking to; in this class, $0 \in \mathbf{N}$.)

Sometime towards the start of junior high you were introduced to the concept of natural numbers. The set

$$\mathbf{Z} = \{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\}$$

is called the set of *integers* or *whole numbers*. (The word *integer* just means “whole” in Latin; cf. French *entier*. We use the letter \mathbf{Z} because of the German word *Zahl* meaning “number.”)

Even before you learned about negative numbers, you probably learned about fractions. They can be defined in set-builder notation as a collection of ratios of integers, where the denominator is not zero:

$$\mathbf{Q} = \{p/q : p, q \in \mathbf{Z}, q \neq 0\}.$$

(The nitpicky reader will notice that p/q is not stipulated to be a member of any set here. This is because a rigorous definition of \mathbf{Q} involves quotienting out by an equivalence relation, which we don't know how to do yet.) This is the set of *rational numbers*. Remember that sets only contain distinct elements, so $2/3$, $4/6$, $6/9$, etc. are all considered the same rational number.

Lastly, we have the set \mathbf{R} of real numbers. Constructing this set using only notions from set theory and logic is quite the byzantine task and well outside the scope of this course, but you can think of \mathbf{R} as the set of decimal numbers with a finite number of digits to the left of the decimal point, and a possibly infinite number of digits to the right of the decimal point.

A set of numbers near and dear to many mathematicians' hearts is the set of *prime* numbers P , defined by

$$P = \{p \in \mathbf{N} : p \geq 2, \text{ and if } p = ab \text{ then } \{a, b\} = \{1, p\}\}.$$

(You may have met a different definition of prime numbers in the past. Pause a moment and convince yourself that the statement above defines the set of prime numbers as you know it.)

Set inclusion. When we use set builder notation $B = \{x \in A : P(x)\}$, every element of B is necessarily a member of the set A as well, since B is defined to be the set of all x in A satisfying $P(x)$. This is one way in which we can obtain a *subset* of another set. More generally, we write $B \subseteq A$ if every element of B is also an element of A , and $B \supseteq A$ if every element of A is an element of B . Sometimes, if $B \supseteq A$, we say that B *contains* A , or B *includes* A . For example, we have the chain

$$\mathbf{N} \subseteq \mathbf{Z} \subseteq \mathbf{Q} \subseteq \mathbf{R}$$

for the special sets of numbers defined earlier.

Symbols like \subseteq , \supseteq , and $=$ (that are used to produce statements) can be negated with a slash; for example, if there is some element of B that is not an element of A , then we write $B \not\subseteq A$.

The concept of set inclusion is important, because the most common way to prove that two sets A and B are equal is to show that A is a subset of B , then show that B is a subset of A . We illustrate this with the following example.

Proposition 1. *The sets*

$$A = \{x \in \mathbf{Z} : \text{there exists } k \in \mathbf{Z} \text{ such that } x = 2k + 1\}$$

and

$$B = \{x \in \mathbf{Z} : \text{there exists } l \in \mathbf{Z} \text{ such that } x = 2l + 5\}$$

are equal. (Both are different ways of expressing the set of all odd integers.)

Proof. Let $x \in A$. Then there exists $k \in \mathbf{Z}$ such that $x = 2k + 1$. Letting $l = k - 2$, we find that l is an integer (since k was). Furthermore,

$$x = 2k + 1 = 2k - 4 + 5 = 2(k - 2) + 5 = 2l + 5.$$

We have found l such that $x = 2l + 5$, so $x \in B$. This shows that $A \subseteq B$.

On the other hand, let $x \in B$, so that there exists $l \in \mathbf{Z}$ such that $x = 2l + 5$. Now we let $k = l + 2$; $k \in \mathbf{Z}$ since $l \in \mathbf{Z}$. We have

$$x = 2l + 5 = 2l + 4 + 1 = 2(l + 2) + 1 = 2k + 1.$$

This shows that $x \in A$, and we have proved that $B \subseteq A$. This combined with the previous paragraph shows that $A = B$. ■

The above result is not so important, but pay attention to the structure of this proof. It is called a “proof by double inclusion,” since we have shown that A includes B and B includes A .

Set operations. Now we describe a number of operations that may be performed on sets to produce other sets. They can all be built up from the following two operations.

- The *union* $A \cup B$ of two sets A and B is the set of all elements that are either in A or in B (or both).
- The *intersection* $A \cap B$ of A and B is the set of all elements that are in both A and B .

As an example, if $A = \{1, 2, 4\}$ and $B = \{1, 3, 5\}$, then $A \cup B = \{1, 2, 3, 4, 5\}$ and $A \cap B = \{1\}$.

To avoid logical difficulties, we always assume that the sets we’re working with are a subset of some larger ambient set U , often called the *universe*. Once we know what U is, we may define the *complement* of a set A to be the set \overline{A} of

all the elements in U except those that are in A . So if $U = \{1, 2, 3, 4, 5\}$ in the example above, then $\overline{A} = \{3, 5\}$ and $\overline{B} = \{2, 4\}$. What about $\overline{A \cup B}$? Well since $A \cup B$ is all of U , its complement must be empty, and we can denote it $\{\}$. This is one valid notation for the *empty set*. The other is \emptyset .

Now is a good time to define the cardinality $|A|$ of a set A . This is the number of elements in it, so $|A| = |B| = 3$ in our example, and $|A \cup B| = 5$, etc. We have $|\emptyset| = 0$, and it is possible for the cardinality of a set to be infinity; for example, $|\mathbf{N}| = \infty$. We also have $|\mathbf{R}| = \infty$, but this infinity is, in some sense, larger than $|\mathbf{N}|$. (More on that later.)

Next, we define the *difference* $B \setminus A$ (sometimes $B - A$) of two sets. This is the set of all elements in B that are *not* in A . So, using the complement notation we just learned about, we can express $B \setminus A = B \cap \overline{A}$. It is not necessary that A be a subset of B . In the small example above, we have $A \setminus B = \{2, 4\}$ and $B \setminus A = \{3, 5\}$.

Lastly, we define the *symmetric difference* $A \triangle B$ of two sets A and B to be the set of all elements that are either in A or in B *but not both*. Invoking the above example one last time, we have $A \triangle B = \{2, 3, 4, 5\}$. We can express it as using unions, intersections, and complements as

$$A \triangle B = (A \cup B) \cap \overline{A \cap B}. \quad (1)$$

To practise using all the different operations we just learned, convince yourself that the following are three more valid ways to express the symmetric difference:

$$A \triangle B = (A \cup B) \setminus (A \cap B) = (A \setminus B) \cup (B \setminus A) = \overline{\overline{A \cup B} \cap (A \cap B)}$$

More set identities abound. We state the following proposition without proof; you should try going through this list and convincing yourself that each identity holds, for all sets A , B , and C . (This is a great way of practising proofs by double inclusion.)

Proposition 2. *Let A , B , and C be subsets of a universe U . Then*

- i) $A \cap U = A$ and $A \cup \emptyset = A$;
- ii) $A \cup U = U$ and $A \cap \emptyset = \emptyset$;
- iii) $A \cup A = A$ and $A \cap A = A$;
- iv) $\overline{\overline{A}} = A$;
- v) $A \cup B = B \cup A$ and $A \cap B = B \cap A$;
- vi) $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$ and $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$;
- vii) $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$ and $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$;
- viii) $A \cup (A \cap B) = A$ and $A \cap (A \cup B) = A$; and
- ix) $A \cup \overline{A} = U$ and $A \cap \overline{A} = \emptyset$. ■

These laws have names, some of which we will use more often than others: (i) is called the identity law, (ii) the domination law, (iii) the idempotent law,

(iv) the law of double negation, (v) the commutative law, (vi) the associative law, (vii) the distributive law, (viii) the absorption law, and (ix) the complement law.

***Analogy with addition and multiplication.** Some of these laws bear some resemblance to laws about numbers that you already know. As an exercise, replace \cup with $+$ (addition), \cap with \cdot (multiplication), \overline{A} with $-A$ (negation), U with 1 , and \emptyset with 0 in all the formulas above, and now assume that A , B , and C are arbitrary real numbers. Which identities still hold in the number setting, and which ones don't? As a more advanced exercise, try replacing \cup with Δ in the identities above (some statements will have to be tweaked a bit so that they're actually true, some won't). Now do the same replacement as before, except replace Δ with $+$. You will find that many more identities carry over.

03.IX **De Morgan's laws.** There are two important laws relating complements with union and intersection. We shall state them as a proposition, this time giving a proof (of one of them).

Proposition 3 (*De Morgan's laws*). *Let A and B be sets. Then*

- i) $\overline{A \cup B} = \overline{A} \cap \overline{B}$; and
- ii) $\overline{A \cap B} = \overline{A} \cup \overline{B}$.

Proof. Let $x \in \overline{A \cup B}$. This means that x does not belong to the union of A and B , x cannot be in A , nor can it be in B . Since $x \notin A$, $x \in \overline{A}$, and since $x \notin B$, $x \in \overline{B}$. Therefore, $x \in \overline{A} \cap \overline{B}$. This shows that $\overline{A \cup B} \subseteq \overline{A} \cap \overline{B}$.

Now assume that $x \in \overline{A} \cap \overline{B}$. So $x \in \overline{A}$ and $x \in \overline{B}$, meaning that $x \notin A$ and $x \notin B$. Since x is in neither A nor B , it is also not a member of the union $A \cup B$. We conclude that $x \in \overline{A \cup B}$. We have shown that $\overline{A} \cap \overline{B} \subseteq \overline{A \cup B}$, which fact, combined with the previous paragraph, completes the proof of (i).

The proof of (ii) is similar and left to the reader as an exercise. ■

Armed with all of these laws, we are able to perform lots of mechanical set manipulations to simplify expressions. For example, consider the expression

$$((A \setminus B) \cup A) \cap \overline{A \cap B}.$$

Since $A \setminus B = A \cap \overline{B}$ and invoking the second De Morgan law on the right of the intersection yields

$$((A \cap \overline{B}) \cup A) \cap (\overline{A} \cup B).$$

Now, we can use absorption on the left-hand side to obtain

$$A \cap (\overline{A} \cup B),$$

and then distributing gives us

$$(A \cap \overline{A}) \cup (A \cap B) = \emptyset \cup (A \cap B) = A \cap B.$$

We thus see that the nasty expression $((A \setminus B) \cup A) \cap \overline{A \cap B}$ is simply another way of writing $A \cap B$.

The Cartesian product and power set. From the real line \mathbf{R} , we can geometrically construct the Cartesian plane \mathbf{R}^2 by lining up parallel copies of \mathbf{R} , one for each element of the original line and all parallel to the original line. Notationally, \mathbf{R}^2 is the set of all ordered pairs (a, b) , where $a, b \in \mathbf{R}$. Generalising this, for any sets A and B we can define the *Cartesian product* $A \times B$ to be the set

$$A \times B = \{(a, b) : a \in A, b \in B\}.$$

We sometimes write A^2 for $A \times A$, and more generally A^n for the n -fold Cartesian product of A with itself. (This explains the notation \mathbf{R}^2 for the Cartesian plane, and \mathbf{R}^n for the n -dimensional vector space over \mathbf{R} .) Note that $A \times B$ is not equal to $B \times A$ in general.

Proposition 4. *If A and B are finite sets, then $|A \times B| = |A| \cdot |B|$.*

Proof. The set $A \times B$ consists of all ordered pairs (a, b) where $a \in A$ and $b \in B$. There are $|A|$ choices for a , and for a , there are $|B|$ ways to pair it with a b from B . So there are $|A| \cdot |B|$ pairs in total. ■

Now we define the *power set*. For a set A , this is the set of all subsets of A , and is commonly denoted by $\mathcal{P}(A)$ or 2^A . (We will use the latter notation in these notes.) In set-builder notation, we have

$$2^A = \{X \subseteq A : X \subseteq A\}.$$

As an example, if $A = \{1, 2, 3\}$, then

$$2^A = \{\emptyset, \{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\}\}.$$

Note that even though, say, $1 \in A$, we do not have $1 \in 2^A$. We do, however, have $\{1\} \in 2^A$.

Another example is $2^{\mathbf{Z}}$, the set of all subsets of integers. If P is the set of primes, then $P \in 2^{\mathbf{Z}}$. Also in $2^{\mathbf{Z}}$ is the set $S = \{n^2 : n \in \mathbf{Z}\}$ of square numbers.

There is a way of encoding subsets with strings of 0s and 1s. Suppose we have a set

$$\{-3, 1, 7, 19, 23\}.$$

Fixing this order of the elements, for any arbitrary subset of this set, we can associate to it a binary string. Consider the subset $\{1, 19, 23\}$. This corresponds to the binary string 01011: since the first element, -3 , is not in the subset, we write a 0. Then 1 is in the subset, so we write a 1, and so on. This is a reversible process. Given a binary string of length 5, say, 00101, we can reconstruct the subset that corresponds to it. The first two 0s mean that -3 and 1 do not belong to the set, but 7 does, 19 doesn't, and 23 does. So the subset is $\{7, 23\}$. In this way we see that there is a one-to-one correspondence between the elements of 2^A and binary strings of length $|A|$. We'll use this fact in the proof of the following proposition.

Proposition 5. *If A is finite, then $|2^A| = 2^{|A|}$.*

Proof. We just saw that there is a one-to-one correspondence between elements of 2^A and binary strings of length $|A|$. So it suffices to count the number of binary strings of length $|A|$. Well, each digit can be either 0 or 1, and there are $|A|$ digits, so the number of strings is

$$\underbrace{2 \cdot 2 \cdot \dots \cdot 2}_{|A| \text{ times}} = 2^{|A|}. \quad \blacksquare$$

Counterexamples in proofs. We finish this subsection with a little example problem. *Is it true that $2^A \cup 2^B = 2^{A \cup B}$ for all sets A and B ?*

Let's start by trying to prove the statement is true. As usual, we will attempt a double-inclusion proof. Let $X \in 2^A \cup 2^B$. This means X is either a subset of A or it is a subset of B . Either way, X is a subset of $A \cup B$, so $X \in 2^{A \cup B}$. So far so good; we have proved that $2^A \cup 2^B \subseteq 2^{A \cup B}$.

Now we try the other direction. Let $X \in 2^{A \cup B}$. So X is a subset of $A \cup B$. From here we want to say that X must be a subset of A or it must be a subset of B , but is that necessarily true? It is possible that X is contained slightly in A and slightly in B . So we have failed to prove that $2^{A \cup B} \subseteq 2^A \cup 2^B$ in general. But just because we have failed to prove that something is true doesn't mean we have proved it is false!

To actually prove that $2^{A \cup B} \subseteq 2^A \cup 2^B$ doesn't hold in general, we need to find a *counterexample*. That is, we need to construct sets A and B such that the statement is false. In this case, we can let $A = \{1, 2\}$, $B = \{3, 4\}$, so that $A \cup B = \{1, 2, 3, 4\}$. Then the set $\{1, 3\}$ is a subset of $A \cup B$ but is not a subset of A and it is not a subset of B . In other words, $\{1, 3\} \in 2^{A \cup B}$ but $\{1, 3\} \notin 2^A \cup 2^B$, proving that $2^{A \cup B} \not\subseteq 2^A \cup 2^B$.

05.IX **Russell's paradox.** Earlier, we said that the sets we are working with need to be a subset of a universe U , which has already been proved to be a set. We gave lots of ways to make sets out of new sets, such as the union and intersection operations, etc. Starting with the assumption that the empty set \emptyset is a set, it is possible to define the set of natural numbers as follows. We can define $0 = \emptyset$, $1 = \{0\}$, and $2 = \{0, 1\}$, and so on. Now we take the set of all of these, and call this \mathbf{N} . (We can also define addition and multiplication on these set-theoretic "numbers" so that they behave like addition and multiplication do on \mathbf{N} .) From here we can do more funky stuff to define \mathbf{Z} , \mathbf{Q} , and \mathbf{R} , and prove that these are all sets (you can see this in a higher-level course on set theory, if you're interested). So there isn't much of a problem with all the sets we have played with so far; they are all subsets of things that are already known to be sets.

Ungodly things can happen if we don't stick by these rules. An example, due to Bertrand Russell, is the "set"

$$R = \{\text{sets } X : X \notin X\}.$$

In plain English, R is defined to be the set of all sets that contain themselves. This is not a subset of any known thing, so by our criterion above we would not consider it a set. But supposing it is, let us ask ourselves the following question. Does R contain itself? If it does not, then $R \notin R$, so R would be a set that satisfied the condition of R , so $R \in R$. But on the other hand, if $R \in R$, then R violates the condition defining R , so $R \notin R$. Round and round we go in a circle of contradiction.

Such is the price of meddling with “sets” that aren’t subsets of known sets. This also shows that there is no such thing as “the set of all sets.”

2. Propositional logic

A *proposition* is a statement that is true or false. For example “8 is even” is a statement we know to be true, and “8 is prime” is a statement we know to be false. The statement “ n is prime” is not a proposition because its truth or falsity depends on what n is. The statement “ $2^{2^{40}} - 1$ is prime” is a proposition, because it is either true or false (even though you or I might not know which one it is).

A *propositional variable* or a *boolean variable* is a variable which can take either the value 0 or 1, where 0 means “false” and 1 means “true.” Usually we use letters p , q , and r to denote propositional variables. The simplest logical operator is negation, defined by the table

p	$\neg p$
0	1
1	0

This is also called the NOT operator, since if p is true, then $\neg p$ is false, and vice versa. Next is *conjunction*, which has the table

p	q	$p \wedge q$
0	0	0
0	1	0
1	0	0
1	1	1

This is also called the AND operator, because $p \wedge q$ is true if and only if p and q are both true. The OR operator, also called *disjunction*, has the table

p	q	$p \vee q$
0	0	0
0	1	1
1	0	1
1	1	1

We see that $p \vee q$ is true if p or q is true (or both). The symbol \vee is meant to recall the Latin word *vel*, meaning “or.” (One of the most important early treatises on mathematical logic and set theory was *Arithmetices principia, nova methodo exposita*, published in Latin in 1889 by Giuseppe Peano. It established the now-standard axiomatisation of the natural numbers.)

On the other hand, in English, we often use the word “or” to mean an *exclusive* or; that is, either p and q are true but not both. In mathematics, on the other hand, “or” is usually *inclusive*, so both p and q are allowed to hold at the same time. It is possible to express the exclusive disjunction (often called XOR) by a table, however. We will use the symbol \oplus for this operator, and its table looks like this:

p	q	$p \oplus q$
0	0	0
0	1	1
1	0	1
1	1	0

So $p \oplus q$ is true if p is true or q is true, but not both.

A *formula* is an expression containing propositional variables, 0, 1, logical operators, and parentheses. The formula must syntactically make sense; for instance, $0(\vee \wedge q \neg$ is not a formula. Just as in ordinary mathematical notation, parentheses are used to clarify which operators should be evaluated first. We will assume that negation applies first, but an expression such as $p \vee q \wedge r$ is ambiguous. (In many programming languages, conjunction has higher priority than disjunction, but in this class, just add parentheses to clarify.)

Above we have illustrated the basic logical operators by writing out their *truth tables*. These are tables that give the value of a formula for all possible values of its variables. We can write truth tables for more complex formulas as well:

p	q	$p \wedge q$	$\neg(p \wedge q)$	$\neg(p \wedge q) \oplus q$
0	0	0	1	1
0	1	0	1	0
1	0	0	1	1
1	1	1	0	1

Strictly speaking, the third and fourth columns are not necessary, but these intermediary columns help us verify the accuracy of the following ones.

Two formulas f_1 and f_2 are said to be *logically equivalent* if they have the same truth table; that is, they produce the same output if given the same input. In this case we write $f_1 \equiv f_2$. For example, let $f_1 = \neg(p \wedge q) \oplus q$, the formula

whose truth table is illustrated above. Now let $f_2 = p \vee \neg q$. Its truth table is

p	q	$\neg q$	$p \vee \neg q$
0	0	1	1
0	1	0	0
1	0	1	1
1	1	0	1

so we conclude that $f_1 \equiv f_2$. As a larger example, suppose we want to find all values of p , q , and r such that

$$f = (p \vee q) \wedge (\neg q \vee \neg r)$$

evaluates to 1. The truth table

p	q	r	$p \vee q$	$\neg q \vee \neg r$	f
0	0	0	0	1	0
0	0	1	0	1	0
0	1	0	1	1	1
0	1	1	1	0	0
1	0	0	1	1	1
1	0	1	1	1	1
1	1	0	1	1	1
1	1	1	1	0	0

shows that f evaluates to 1 precisely when

$$(p, q, r) \in \{(0, 1, 0), (1, 0, 0), (1, 0, 1), (1, 1, 0)\}.$$

To write out the truth table for a formula with n variables, we need 2^n rows, so this method is unsuitable for formulas with more than three or four variables.

Simplifying logical formulas. Just as we have rules for simplifying set expressions, there are ways to turn complicated logical formulas into simpler ones that are logically equivalent. What might surprise you is that the rules turn out to be exactly the same! To see the basis for this correspondence, consider the definition of a union of sets A and B . In set-builder notation, this is

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

The “or” in the definition suggests that \cup is intimately related to the \vee operation in propositional logic. Repeating this process, we have the “dictionary”

Set theory	Propositional logic
sets A, B	variables p, q
unions $A \cup B$	disjunctions $p \vee q$
intersections $A \cap B$	conjunctions $p \wedge q$
complements \overline{A}	negations $\neg p$
symmetric differences $A \triangle B$	exclusive disjunctions $p \oplus q$
the empty set \emptyset	0
the universe U	1

Exploiting this connection, we have the following analogue of Proposition 2.

Proposition 2. *Let p , q , and r be propositional variables. Then*

- i) $p \wedge 1 \equiv p$ and $p \vee 0 \equiv p$;
- ii) $p \vee 1 \equiv 1$ and $p \wedge 0 \equiv 0$;
- iii) $p \vee p \equiv p$ and $p \wedge p \equiv p$;
- iv) $\neg\neg p \equiv p$;
- v) $p \vee q \equiv q \vee p$ and $p \wedge q \equiv q \wedge p$;
- vi) $p \vee (q \vee r) \equiv (p \vee q) \vee r$ and $p \wedge (q \wedge r) \equiv (p \wedge q) \wedge r$;
- vii) $p \vee (q \wedge r) \equiv (p \vee q) \wedge (p \vee r)$ and $p \wedge (q \vee r) \equiv (p \wedge q) \vee (p \wedge r)$;
- viii) $p \vee (p \wedge q) \equiv p$ and $p \wedge (p \vee q) \equiv p$; and
- ix) $p \vee \bar{p} \equiv 1$ and $p \wedge \bar{p} \equiv 0$. ■

Since they are essentially the same as their set equivalents, the names of these laws are the same as in the realm of sets. We also have the propositional equivalent of De Morgan's law, which states that

$$x) \neg(p \vee q) \equiv \neg p \wedge \neg q \text{ and } \neg(p \wedge q) \equiv \neg p \vee \neg q.$$

Using these rules, we can now show that $\neg(p \wedge q) \oplus q \equiv p \vee \neg q$, which we already saw earlier from their truth tables. First we observe that

$$p \oplus q \equiv (p \vee q) \wedge \neg(p \wedge q), \quad (2)$$

which is analogous to the identity (1) for symmetric differences. So

$$\begin{aligned} \neg(p \wedge q) \oplus q &\equiv (\neg(p \wedge q) \vee q) \wedge \neg(\neg(p \wedge q) \wedge q) \\ &\equiv (\neg p \vee \neg q \vee q) \wedge ((p \wedge q) \vee \neg q) \\ &\equiv (\neg p \vee 1) \wedge ((p \vee \neg q) \wedge (q \vee \neg q)) \\ &\equiv 1 \wedge ((p \vee \neg q) \wedge 1) \\ &\equiv p \vee \neg q, \end{aligned}$$

where in the first line we use (2), in the second line we use De Morgan's law twice, in the third line we use the complement and distributive laws, the fourth line we use the domination and complement laws, and in the last line we use the identity law (twice). (It is not super important to remember the names of all these laws, as long as you remember their statements, but you may actually find it is easier to remember the names along with the statements, just as it might be easier to remember faces of people you've met if you also know their names.)

10.IX **Conditional and biconditional.** We now examine the *conditional* logical relation IF p THEN q . It is denoted by $p \Rightarrow q$ and its truth table is given by

p	q	$p \Rightarrow q$
0	0	1
0	1	1
1	0	0
1	1	1

Within the conditional statement, p is called the *antecedent*; this is the assumption. The statement that is asserted, conditional that the antecedent holds, is called the *consequent* q . You can quickly check that the relation $p \Rightarrow q$ is equivalent to $\neg p \vee q$. This fact is useful when performing mechanical simplifications. In English, the statement “if p then q ” asserts a causal relation between p and q . Take a moment and reconcile this idea with the truth table above. It is normal to get a little tripped up if it’s your first time seeing this table. In the first row, q doesn’t even happen, so it might feel weird to say that p has “caused” q to happen in this case, and in the second row, p is not true and q is true, so it might seem odd that we have set $p \Rightarrow q$ to true, because there seems to be no relation between p and q .

But it should make sense if you think of p as a precondition to a promise q , and then consider whether the promise is broken. As an example, suppose your friend says, If it snows tomorrow I’ll work in Trottier with you. If it doesn’t snow and she doesn’t pull up, she hasn’t technically broken her promise. If it doesn’t snow and she shows up, then she still hasn’t broken her promise. The only way she can break her promise is if it snows and she doesn’t come to Trottier; this situation corresponds to the only 0-row in the truth table.

The conditional is a very important logical operator to understand, because most theorem statements assume some hypothesis and claim some conclusion. You will be asked to prove statements of this form, so it is important to understand the logical nature of the statements to begin with.

The last relation is called the *biconditional*, and it asserts that variables p and q are logically equivalent. That is, p happens if and only if q happens. It’s truth table

p	q	$p \Leftrightarrow q$
0	0	1
0	1	0
1	0	0
1	1	1

is pretty self-explanatory; in it, $p \Leftrightarrow q$ is true whenever p and q have the same truth value. The name “biconditional” is suggested by (part of) the following proposition.

Proposition 6. *Let p and q be propositional variables. Then*

$$\begin{aligned}
 p \Leftrightarrow q &\equiv (p \Rightarrow q) \wedge (q \Rightarrow p) \\
 &\equiv ((\neg p) \vee q) \wedge ((\neg q) \vee p) \\
 &\equiv (p \wedge q) \vee ((\neg p) \wedge (\neg q)).
 \end{aligned}$$

Proof. We leave the first equivalence to the reader (writing out the truth table is one way of proving it). The second equivalence follows from our earlier observation that $p \Rightarrow q \equiv \neg p \vee q$, and the third equivalence follows from the distributive, complement, and identity laws. ■

A formula f is called a

- i) *tautology* if $f \equiv 1$, i.e., f always evaluates to true;
- ii) *contradiction* if $f \equiv 0$, i.e., f always evaluates to false;
- iii) *contingency* if f can evaluate to both 1 and 0, depending on the values of its variables;
- iv) *satisfiable* if f evaluates to 1 for at least one input; and
- v) *falsifiable* if f evaluates to 0 for at least one input.

An example of a tautology is $p \vee \neg p$ and an example of a contradiction is $p \wedge \neg p$. This follows from the complement laws. To say that something is satisfiable is precisely to say that it is not a contradiction, and to say that something is falsifiable is equivalent to saying that it is not a tautology. Contingencies are those formulas that are both satisfiable and falsifiable (in other words, formulas that are neither tautologies nor contradictions).

Suppose we are asked which of the above definitions the formula

$$f \equiv (p \wedge (p \Rightarrow q)) \Rightarrow q$$

satisfies. This formula only has two variables, so it is easy enough to use a truth table for this purpose, but we will take the opportunity to practise simplifying the expression symbolically. First of all, let's change all conditionals of the form $r \Rightarrow s$, to disjunctions of the form $\neg r \vee s$. This gives us

$$f \equiv \neg(p \wedge (\neg p \vee q)) \vee q.$$

Now we use the distributive law to distribute the conjunction over the innermost disjunction, obtaining

$$f \equiv \neg((p \wedge \neg p) \vee (p \wedge q)) \vee q;$$

by the complement and identity laws in that order, this simplifies to

$$f \equiv \neg(p \wedge q) \vee q.$$

Now De Morgan's law and associativity give

$$f \equiv (\neg p \vee \neg q) \vee q \equiv \neg p \vee (\neg q \vee q),$$

and thus

$$f \equiv \neg p \vee 1 \equiv 1,$$

by the complement and domination laws in that order. We conclude that f is a tautology, which also means that it is satisfiable.

The fact that $(p \wedge (p \Rightarrow q)) \Rightarrow q$ is a tautology symbolically justifies the argument that if p is true and $p \Rightarrow q$ is true, then we should be able to conclude q . This form of argument is called *modus ponens*, and it dates back to ancient

times. You probably use *modus ponens* all the time in everyday life without knowing it, and we will certainly use it in this class a lot.

Encoding problems in propositional logic. Many algorithmic and logical problems can be encoded in propositional logic (and then later solved by a computer program). For example, suppose we want to play 4×4 Sudoku. In this game, we have a 4×4 grid and we want to fill it with the numbers 1 through 4 such that

- i) every row contains 1 through 4;
- ii) every column contains 1 through 4; and
- iii) the four subsquares each contain 1 through 4.

In a given instance of the game, some cells are already filled in. The puzzle is: *Is there a solution and if so, what is it?*

		2	
1	3		

Fig. 1. An example 4×4 Sudoku game.

To represent a Sudoku game in propositional logic, we can define boolean variables $p_{i,j,k}$, where i , j , and k range over $\{1, 2, 3, 4\}$. (So there are 4^3 variables in total.) We shall set

$$p_{i,j,k} = \begin{cases} 1, & \text{if the number } k \text{ is in row } i \text{ and column } j; \\ 0, & \text{otherwise.} \end{cases}$$

We'll number the rows increasing from the top and the columns increasing from left to right. For example, in Fig. 1 there is a 2 in row 1 and column 3, so $p_{1,3,2} = 1$.

Now we set to work encoding the conditions of a Sudoku grid in propositional logic:

- i) To stipulate that every row contain 1 through 4, we first define auxiliary variables

$$r_{i,k} = p_{i,1,k} \vee p_{i,2,k} \vee p_{i,3,k} \vee p_{i,4,k},$$

for $i, k \in \{1, 2, 3, 4\}$. With these helper variables, we now see that

$$r_{1,1} \wedge r_{1,2} \wedge r_{1,3} \wedge r_{1,4}$$

encodes the requirement that row 1 contains one of each number. We do the same for rows 2, 3, and 4 as well, and then combine with AND.

- ii) We do a similar thing as in part (i) for each of the four columns.
- iii) Ditto for subsquares.
- iv) We need to set the initial values of the grid. For the grid in Fig. 1, we have the formula

$$p_{1,3,2} \wedge p_{3,1,1} \wedge p_{3,2,3}.$$

- v) Lastly, we need to make sure that there is not more than one number per cell. To do this, for each cell $(i, j) \in \{1, 2, 3, 4\}^2$ we write

$$p_{i,j,1} \Rightarrow (\neg p_{i,j,2} \wedge \neg p_{i,j,3} \wedge \neg p_{i,j,4}),$$

and so on (four conditionals in total). Of course we'll need to AND all these together.

Now we combine the formulas from each of these five steps into one long formula f such that f is satisfiable if and only if the grid has a solution, and the values for $p_{i,j,k}$ give a solution.

Defining all of these variables was a rather arduous and cumbersome process, and not entirely worth it for a 4×4 game of Sudoku (which can just be solved by eyeballing the grid). But one could imagine writing a general computer program to encode larger and larger grids. In fact, there are lots of problems that can be *reduced* to the problem of determining if a boolean formula is satisfiable. This means that if we have a program capable of taking a formula f as input and spitting out whether or not it is satisfiable (in a reasonable amount of time), then there are lots of real-world problems that this program could be applied to.

This problem is called the *boolean satisfiability problem*, often abbreviated SAT. One way of solving it for any given f is to just compute its truth table. We already know the downside of this approach: if f has n variables, then its truth table will have 2^n rows. Given a few minutes, you are certainly capable of writing down a formula f that has 300 variables, call them p_1, \dots, p_{300} . The truth table of f will have 2^{300} rows, which is more than the number of atoms in the observable universe. You can learn a lot more about SAT in a higher-level class on computational complexity (e.g., COMP 360/362).

3. Predicate logic

- 12.IX Propositional logic allows us to work with simple declarations, but this isn't powerful enough to express some deeper mathematical concepts. For this purpose, we now introduce the notion of a *predicate*. This is a statement involving some number of variables, each of which may take values coming from a universe U , such that the statement evaluates to either true or false *once all variables are assigned values*. The statement $P(n)$ given by " n is prime" is an example of this, where n can take any value in the universe \mathbf{Z} . An example with two variables is the predicate $L(x, y)$ defined by " x is less than y ." (A more commonly-used notation for this predicate is " $x < y$.")

Predicates contain variables, but at the moment we don't have any way of introducing new variables into a statement. This is done using two different *quantifiers*. The first is the *universal quantifier*, denoted \forall and meaning “for all.” The statement $\forall n : P(n)$ is true if and only if $P(n)$ is true for every possible value that n can take. The second quantifier is the *existential quantifier*, written “ \exists ” and with the meaning “there exists.” The statement $\exists n : P(n)$ is true if and only if there is (at least) one value that n can take such that $P(n)$ is true. The colon doesn't really have any mathematical meaning in these formulas; they just visually set the quantifiers apart from the predicates that follow.

For instance, taking $P(n)$ to be the statement “ n is prime,” where the universe U is \mathbf{N} , the statement $\exists n P(n)$ is true and the statement $\forall n : P(n)$ is false. What about the statement $\forall x \exists y : y < x$? Well, if the universe U is taken to be \mathbf{N} , then the statement is false, because setting x equal to 0, there is no element y of \mathbf{N} such that $y < 0$. But if $U = \mathbf{Z}$, then the statement is true, since for every integer x , we can put $y = x - 1$, so that $y < x$.

Now we practise translating converting mathematical statements written in English into formulas in predicate logic. Suppose we want to write, Every integer is even or odd. The universe here is the set \mathbf{Z} of integers. The word “every” has the same meaning as “for all,” so right off the bat, we can reexpress the statement as, For all $n \in \mathbf{Z}$, n is even or n is odd. In symbols, this is

$$\forall n (n \text{ even} \vee n \text{ odd}).$$

Lastly, we need to figure out how to express the property of being even or being odd. An integer n is even if and only if it is a multiple of two; that is, if there is some integer k such that $2k = n$. Likewise, an integer is odd if and only if it is one more than a multiple of two. The corresponding formula is $\exists k : 2k + 1 = n$. So our statement can be expressed

$$\forall n ((\exists k : n = 2k) \vee (\exists k : n = 2k + 1)).$$

The variable k appears twice in this formula, but its first instance is independent of its second instance, because of the parentheses. (Readers who write computer programs will be familiar with the concept of a variable “going out of scope.”) So there is nothing wrong with this formula, but to be extra clear that the first k is different from the second k , why don't we replace it with a different letter? Thus we arrive at

$$\forall n ((\exists k : n = 2k) \vee (\exists l : n = 2l + 1)),$$

a formula in predicate logic that means, Every integer is even or odd.

Restrictions using quantifiers. It is not true that every real number has a multiplicative inverse, since one cannot divide by zero. However, the statement “every nonzero real number has a multiplicative inverse” is true. How should

we write this as a formula in predicate logic? One way is to use the conditional: over the universe $U = \mathbf{R}$, we could write

$$\forall x : (x \neq 0 \Rightarrow \exists y : xy = 1).$$

Another is to use subscripts: in the same universe, we write

$$\forall x_{x \neq 0} \exists y : xy = 1.$$

Using subscripts is slightly informal, since we didn't formally define above what a subscript is supposed to mean, but it is something that you might encounter. The last way is to simply restrict the universe itself: in the universe $U = \mathbf{R} \setminus \{0\}$, the statement

$$\forall x \exists y : xy = 1$$

is true.

Multiple quantifiers. Withing a formula, quantifiers cannot be interchanged willy-nilly. The order of \forall and \exists matters! They are read from left to right. Consider the following examples, over the universe $U = \mathbf{R}$. The statement

$$\forall x \exists y : x + y = 0$$

is true, since for each given x we can take y to be $-x$. On the other hand,

$$\exists y \forall x : x + y = 0$$

is false, since it would mean that there is some integer that adds up to zero with *any* integer. Of course, sometimes switching the order of quantifiers, doesn't change the truth value of a statement. Both

$$\exists y \forall x : xy = 0$$

and

$$\forall x \exists y : xy = 0$$

are true, since in the first case, we can take $y = 0$, and in the second case, we can set y to 0 no matter what x is given.

So we know that the order of \forall and \exists matters in general, but repeated instances of the *same* quantifier *can* be interchanged. For instance,

$$\forall x \forall y : x^2 + y^4 \geq 0$$

is the same as

$$\forall y \forall x : x^2 + y^4 \geq 0,$$

and we can even write $\forall x, y : x^2 + y^4 \geq 0$, to introduce both variables simultaneously.

Negating quantifiers. The universal quantifier is sort of like a big chain of conjunctions that goes over the whole of the universe. For example, in the universe \mathbf{N} , the statement $\forall n : P(n)$ is equivalent to

$$P(1) \wedge P(2) \wedge P(3) \wedge \cdots,$$

if this were a valid propositional formula (it isn't because we don't allow propositional formulas to be infinite). Likewise, the existential quantifier $\exists n : P(n)$ is equivalent to

$$P(1) \vee P(2) \vee P(3) \vee \cdots.$$

We know, by De Morgan's laws, that negating a big series of conjunctions requires us to flip all the ANDs to ORs. So, once again abusing notation somewhat, we expect

$$\neg(P(1) \wedge P(2) \wedge P(3) \wedge \cdots) \equiv \neg P(1) \vee \neg P(2) \vee \neg P(3) \vee \cdots.$$

Thus we conclude that

$$\neg(\forall n : P(n)) \equiv \exists n : \neg P(n).$$

You can play the same game with the other De Morgan's law to show that

$$\neg(\exists n : P(n)) \equiv \forall n : \neg P(n).$$

Going back to our example of $P(n)$ denoting “ n is prime,” the statement $\neg(\forall n : P(n))$ is true, since not all integers n are prime, and we have just shown that this is equivalent to the statement $\exists n : \neg P(n)$; that is, there exists n such that n is not prime.

We end this section with a longer example. Let's express the statement, “There is a nonzero real number such that every real number is not its inverse or is negative.” In the universe \mathbf{R} , the formula corresponding to this statement is

$$\exists x : (x \neq 0 \wedge (\forall y : xy \neq 1 \vee y < 0)).$$

(Work it out yourself!) Is this statement true or false? It turns out that it is true. You might be able to stare at the formula long enough to convince yourself of this fact, but another way to see that it's true is to note that its negation is false. Let's do this now (it's a good excuse to practise negating a formula). We have

$$\begin{aligned} & \neg(\exists x : (x \neq 0 \wedge (\forall y : xy \neq 1 \vee y < 0))) \\ & \equiv \forall x : \neg(x \neq 0 \wedge (\forall y : xy \neq 1 \vee y < 0)) \\ & \equiv \forall x : (x = 0 \vee \neg(\forall y : xy \neq 1 \vee y < 0)) \\ & \equiv \forall x : (x = 0 \vee \exists y : \neg(xy \neq 1 \vee y < 0)) \\ & \equiv \forall x : (x = 0 \vee \exists y : (xy = 1 \wedge y \geq 0)). \end{aligned}$$

This negated statement is false, since if $x = -2$, then $x = 0$ doesn't hold, so the left-hand side of the OR isn't true, and there is no y such that $-2y = 1$ and $y \geq 0$ are both true, since the only y satisfying $-2y = 1$ is $-1/2$.

Negating a formula in predicate logic is entirely mechanical. The \neg symbol moves from left to right like a bulldozer that flips quantifiers and negates predicates it finds along the way, until eventually its job is done and it disappears. More broadly, we write mathematical statements in formal logic to make things more precise and mechanical. This can be useful to humans, since English is often ambiguous whereas the notation we just established is not. It can be useful to machines as well, since, as we just saw, manipulating a formula is something that can very easily be automated.

4. Proofs

We're getting to the fun part of the course now. Further back in these notes, we already proved a few statements about sets. This was just a taste of what's to come, as the main focus of this course is to teach you how to prove mathematical statements. These will all be statements that can ultimately be stated in predicate logic, so using the tools from the previous section, you can boil them down to their logical skeleton. In this section, we will formally define what it means to prove a statement in predicate logic.

To prove a statement, we always process the quantifiers from left to right. Whenever we encounter something of the form $\exists x : P(x)$, we are allowed to choose the value for x (from the given universe), and we just need to show that $P(x)$ holds for that value of x . Here's an example.

Proposition 7. *There exists an integer $m > 0$ and an integer $n < 0$ such that $m^2 + n^2 = 25$.*

We've written the statement in words, and we will continue to do so for all the propositions in these notes because we are humans and not cyborgs, but notice that the underlying predicate logical formula here is

$$\exists m \exists n : m > 0 \wedge n < 0 \wedge m^2 + n^2 = 25,$$

with $U = \mathbf{Z}$. For the remainder of this section, we'll continue to write out the formulaic equivalents of propositions, to practise converting between the two worlds.

Proof. After a moment's contemplation, we notice that $9 + 16 = 25$. So we need a positive integer that squares to 9 and a negative number that squares to 16. Hence we may pick $m = 3$ and $n = -4$, and $m^2 + n^2 = 25$. ■

On the other hand, to prove a statement of the form $\forall x : P(x)$, we are not allowed to choose the value of x . Instead, we imagine it is given to us, and we still have to prove $P(x)$, no matter what the x might be. Here's what we mean.

Proposition 8. *For all $x \in \mathbf{Q}$, there exists $y \in \mathbf{Z}$ such that $xy \in \mathbf{Z}$.*

The formula this time is

$$\forall x \exists y : y \in \mathbf{Z} \wedge xy \in \mathbf{Z},$$

over $U = \mathbf{Q}$.

Proof. Let $x \in \mathbf{Q}$. Then x can be expressed as the ratio of two integers; write $x = m/n$ with $m, n \in \mathbf{Z}$. Then, setting $y = n \in \mathbf{Z}$, we have $xy = (m/n) \cdot n = m \in \mathbf{Z}$. **■**

Notice how we introduced the variable x in the above proof. Because we're trying to prove a "for all" statement, we use the word "let," to indicate that x is given to us by some higher power. In fact, we should sort of view this higher power as possibly being malicious. A very useful way to think about writing a proof is to imagine a game in which you are trying to prove a statement in predicate logic, and a supernatural adversary is attempting to thwart you. Let's say this predicate has four variables, so the statement is

$$\exists x \forall y \forall z \exists w : P(x, y, z, w).$$

The variables in the statement are introduced from left to right. Each time you see a \exists symbol, it's your turn. In the example above, you get to set the variable x to any element of the given universe (keeping in mind that your eventual goal is to prove $P(x, y, z, w)$). Each time there is a \forall symbol, it's the adversary's turn, and you should be prepared for whatever he throws at you. So in our example, the adversary may set y and z to anything in the universe, and he knows you picked x . Lastly, you get to pick w . If $P(x, y, z, w)$ is true, you win, and if not, the adversary wins. Writing a proof is equivalent to describing a winning strategy for the player against the adversary.

17.IX **Proving conditional statements.** Many mathematical statements introduce some hypotheses, then assert some conclusion. Thus they are some kind of statement of the form $p \Rightarrow q$. To prove this kind of statement, we assume that p holds, then prove that q is true. This is because the only way a statement for $p \Rightarrow q$ to be false is if p is true and q is false, so we're showing this cannot happen. Take a look at the following example.

Proposition 9. *If n is an odd integer, then n^2 is also odd.*

The underlying formula is $\forall n((\exists k : n = 2k + 1) \Rightarrow (\exists l : n^2 = 2l + 1))$, in the universe $U = \mathbf{Z}$.

Proof. Let n be an arbitrary odd integer, so that there exists k such that $n = 2k + 1$. Then

$$n^2 = (2k + 1)^2 = 4k^2 + 4k + 1 = 2(2k^2 + 2k) + 1,$$

so setting $l = 2k^2 + 2k$, which is an integer, we have $n^2 = 2l + 1$. Hence n^2 is odd. ■

To disprove a statement, we simply prove that its negation is true. For example, suppose we want to disprove the statement $\exists x \forall y : x + y = 0$ (this was a statement we encountered in the previous section). Its negation is

$$\forall x \exists y : x + y \neq 0,$$

and to prove it we let x be given, set $y = -x + 1$, and observe that $x + y = x + (-x + 1) = 1$, which is not equal to 0.

Now we define two words that sound similar but are very different, logically speaking. The *converse* of a conditional statement $p \Rightarrow q$ is the statement $q \Rightarrow p$. The statement $p \Rightarrow q$ and its converse both hold if and only if the biconditional statement $p \Leftrightarrow q$ holds. On the other hand, the *contrapositive* to $p \Rightarrow q$ is the statement $\neg q \Rightarrow \neg p$. What's the deal with this silly-looking conditional? Well, it turns out to actually be equivalent to $p \Rightarrow q$. Check it out:

$$\begin{aligned} p \Rightarrow q &\equiv \neg p \vee q \\ &\equiv q \vee \neg p \\ &\equiv \neg(\neg q) \vee \neg p \\ &\equiv \neg q \Rightarrow \neg p \end{aligned}$$

So to prove a statement of the form $p \Rightarrow q$, it is sufficient to prove $\neg q \Rightarrow \neg p$. Sometimes, this can make the proof a lot easier. Here's an example. (It is the converse of Proposition 9.)

Proposition 10. *If n is an integer such that n^2 is odd, then n is odd.*

Proof attempt. Assume n^2 is odd, so that $n^2 = 2k + 1$ for some integer k . Then $n = \sqrt{2k + 1}$.

Where do we go from here? It's not even clear that $\sqrt{2k + 1}$ is an integer, let alone an odd one. The contrapositive comes to the rescue.

Proof of Proposition 10. We proceed by contraposition. Suppose that n is not odd; that is, n is even. So there exists an integer k such that $n = 2k$. Then $n^2 = (2k)^2 = 4k^2$. Setting $l = 2k^2$, we have $n^2 = 2l$, so n^2 is even. Hence n^2 is not odd. ■

Proofs by contradiction. We now discuss a powerful method of proof, which dates back to at least the ancient Greeks and referred to for much of Western history by its Latin name, *reductio ad absurdum*. It proceeds to prove a statement p by assuming its negation $\neg p$ holds, then deriving a contradiction, i.e., a statement whose truth value is 0. By doing this, one will have proved that $\neg p \Rightarrow 0$, which in turn shows that p is true, since

$$\neg p \Rightarrow 0 \equiv \neg(\neg p) \vee 0 \equiv p \vee 0 \equiv p.$$

Here is an example of a proof by contradiction.

Proposition 11. *There is no least positive rational number.*

In $U = \mathbf{Q}$, one way to formulate this is

$$\neg(\exists x : x > 0 \wedge (\forall y : y > 0 \Rightarrow x \leq y)).$$

Proof. Suppose, towards a contradiction, that there exists some rational $x > 0$ such that for all $y \in \mathbf{Q}$ with $y > 0$, $x \leq y$. Then, we can apply this property with $y = x/2$ to see that $x \leq x/2$. Since x is positive, we can divide this inequality by x on both sides to obtain $1 \leq 1/2$. But this is absurd, since $1/2 < 1$. This contradiction completes the proof. ■

Here is another classical example of a proof by contradiction. It was known to the Pythagoreans.

Theorem 12. *The number $\sqrt{2}$ is irrational.*

Proof. Suppose, towards a contradiction, that $\sqrt{2}$ is rational. Then there are integers p and q such that $\sqrt{2} = p/q$, and furthermore, we can assume that p and q do not have any common factors, since if they did, we could divide them out and the ratio would remain the same. Squaring both sides of the equation, we have $2 = p^2/q^2$, or in other words, $2q^2 = p^2$. This implies that p^2 is even, so by (the contrapositive of) Proposition 10, p is even as well. That means we can write $p = 2r$ for some integer r , and substitute this new information into the above equation to obtain $2q^2 = (2r)^2 = 4r^2$. Dividing out by 2 yields $q^2 = 2r^2$, so q^2 is even, and so is q .

We have deduced that p and q are both even. On the other hand, we assumed they had no common factors. This is a contradiction. ■

Earlier we described the composition of a proof as playing a game with a supernatural adversary. To prove something by contradiction, then is akin to starting the game by letting adversary believe he has already won, and then working from that assumption to derive an impossibility. Quite the devious stratagem.

Reductio ad absurdum, which Euclid loved so much,
is one of a mathematicians finest weapons.
It is a far finer gambit than any chess gambit:
a chess player may offer the sacrifice of a pawn
or even a piece,
but a mathematician offers the game.

— G. H. HARDY, *A Mathematician's Apology* (1940)

Case analysis. For all $n \geq 2$ and all propositions p_1, \dots, p_n, q , we have the equivalence

$$(p_1 \vee p_2 \vee \dots \vee p_n) \Rightarrow q \equiv (p_1 \Rightarrow q) \wedge (p_2 \Rightarrow q) \wedge \dots \wedge (p_n \Rightarrow q).$$

(As an exercise, prove this when $n = 2$. Later on, when we learn about mathematical induction, you'll be able to prove it in general.)

Why is this useful? Well, if we have two or more propositions p_1, p_2, \dots, p_n such that at least one of them must hold, that is,

$$p_1 \vee p_2 \vee \dots \vee p_n \equiv 1,$$

then if we are able to show that

$$p_1 \Rightarrow q, \quad p_2 \Rightarrow q, \quad \dots, \quad \text{and} \quad p_n \Rightarrow q$$

all hold, then we will have shown

$$1 \Rightarrow q \equiv \neg 1 \vee q \equiv 0 \vee q \equiv q.$$

Let's see an example of this in action.

Proposition 13. For all integers $n \geq 0$, $1 + (-1)^n(2n - 1)$ is a multiple of 4.

Proof. Let an integer $n \geq 0$ be given. We know that n is even or n is odd.

If n is even, then $n = 2k$ for some integer k , and

$$1 + (-1)^{2k}(2(2k) - 1) = 1 + 1^k(4k - 1) = 4k,$$

which is a multiple of 4.

If n is odd, then $n = 2k + 1$ for some integer k , and

$$1 + (-1)^{2k+1}(2(2k + 1) - 1) = 1 - (4k + 2 - 1) = 1 - (4k + 1) = -4k,$$

which is also a multiple of 4. ■

19.IX Sometimes, the cases into which one might split the proof are not so obvious.

Proposition 14. There exist irrational numbers a and b such that a^b is rational.

Proof. Consider the number $(\sqrt{2})^{\sqrt{2}}$. It is either rational or irrational.

If it is rational, then we can set $a = b = \sqrt{2}$, which we already know to be irrational, by Theorem 12.

If it is irrational, then we can set $a = (\sqrt{2})^{\sqrt{2}}$ and $b = \sqrt{2}$. We compute

$$a^b = ((\sqrt{2})^{\sqrt{2}})^{\sqrt{2}} = (\sqrt{2})^{\sqrt{2} \cdot \sqrt{2}} = (\sqrt{2})^2 = 2,$$

which is rational. ■

Note that in this proof, we did not need to know whether $(\sqrt{2})^{\sqrt{2}}$ is rational or not; we just know that it must either be rational or irrational. (In fact, it is known that $(\sqrt{2})^{\sqrt{2}}$ is irrational, but the usual proof of this requires some notions from complex analysis as well as Galois theory, subjects that are typically not encountered until at least the second or third year of a mathematics degree.)

Mathematical induction. Sometimes we want to prove an infinite number of statements, indexed by the natural numbers. If the complexity of the statements sort of “grows” in n (a vague notion to be made precise soon), the following theorem holds.

Theorem 15 (*Principle of mathematical induction*). Let $P(n)$ be a family of predicates indexed by $n \in \mathbf{N}$. Let $m \in \mathbf{N}$. If

- i) $P(m)$ holds; and
- ii) for all $n \geq m$,

$$(P(m) \wedge P(m+1) \wedge \cdots \wedge P(n)) \Rightarrow P(n+1)$$

then $P(n)$ holds for all $n \geq m$.

In the following proof, we use the *well-ordering principle*, which is the fact that every nonempty subset of the natural numbers has a least element. You can see the proof of the well-ordering principle in a higher-level set theory course.

Proof. Assume that (i) and (ii) both hold. Our goal is to prove that $P(n)$ is true for all $n \geq m$, so we shall suppose, towards a contradiction, that there is some $n \geq m$ such that $\neg P(n)$ holds. (As an exercise, convince yourself, via symbolic manipulations, that

$$\neg \forall n : n \geq m \Rightarrow P(n) \equiv \exists n : n \geq m \wedge \neg P(n).)$$

In other words, the set of $n \geq m$ such that $\neg P(n)$ holds is nonempty, so it has a least element, call it k . So $k \geq m$, and we cannot have $k = m$ due to (i), so $k > m$. Furthermore, by the minimality of k , the statements $P(m)$, $P(m+1)$, all the way up to $P(k-1)$ are true. By (ii) though, this implies that $P(k)$ holds: a contradiction. ■

When proving something by induction, we need to prove two things: one of the form (i) and another of the form (ii). The former is called the *base case* and the latter the *induction step* or *inductive step*. The assumption to the left of the \Rightarrow symbol in (ii) is called the *induction hypothesis* or *inductive hypothesis*. Technically, the theorem above is called the principle of *strong* induction; the principle of *weak* induction has

ii)' for all $n \geq m$, $P(n) \Rightarrow P(n+1)$,

instead of the stronger (ii). It turns out that both types of induction are actually equivalent, so we'll use both interchangeably. More often than not, the hypothesis (ii)' is perfectly sufficient, and in these cases we'll simply use weak induction so as not to clutter our proofs with lots of unused hypotheses.

Here is an example. Suppose we want to find (and prove) a formula for the sum of the first n odd numbers. The first thing to do is to try some small cases. When $n = 1$, the sum is just 1. When $n = 2$, the sum is $1 + 3 = 4$, when $n = 3$, we have $1 + 3 + 5 = 9$, and for $n = 4$, we compute $1 + 3 + 5 + 7 = 16$. So the pattern goes 1, 4, 9, 16, ..., which leads us to conjecture that the sum of the first n odd numbers might equal n^2 . In fact, this is true, and we shall prove it by induction.

Proposition 16. *For all integers $n \geq 1$,*

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

Proof. By induction on n . First we prove the base case, $n = 1$. We have

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

Now for the inductive step. Let $n \geq 1$ and assume that

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

Then

$$\begin{aligned} \sum_{i=1}^{n+1} (2i - 1) &= \sum_{i=1}^n (2i - 1) + (2(n+1) - 1) \\ &= n^2 + 2n + 2 - 1 \\ &= n^2 + 2n + 1 \\ &= (n+1)^2, \end{aligned}$$

where it is in the second line that we used the inductive hypothesis. **■**

The first thing to notice about doing a proof by induction is that the proof method itself doesn't tell you what it is you should prove. You have to guess at the correct statement first. Also, proofs by induction are often "unenlightening," in that they often don't reveal the fundamental reasons why something might be true. (The previous proposition can be illustrated by a rather simple picture, which is not a proof, but is somewhat more elucidating than the induction proof.)

A longer example now. Suppose we have a pile of n stones, with $n \geq 1$. We have a job, which can be described by a pseudo-algorithm.

Algorithm S (*Divide stones*). The input to this algorithm is an integer $n \geq 1$, representing a number of stones. We have a list **PILES** of integers representing a collection of piles of stones, as well as an integer variable **BALANCE**. Initialise **PILES** $\leftarrow 1$ (one pile with n stones), and set **BALANCE** $\leftarrow 0$. This algorithm splits the stones into n piles of 1 stone each, accumulating profits into **BALANCE** along the way.

- S1.** [Done?] If every element of **PILES** is 1, terminate the algorithm and output **BALANCE**.
- S2.** [Choose a pile.] Select some element $m > 2$ from **PILES** and remove it from the list. (The variable m is the number of stones in this pile.)

S3. [Split.] Let k and l be two numbers with $k + l = m$. We append $\text{PILES} \leftarrow k$ and $\text{PILES} \leftarrow l$, and increment $\text{BALANCE} \leftarrow \text{BALANCE} + k \cdot l$. (We have split pile n into two piles of size k and l , and the payout for doing so is $k \cdot l$ dollars.)

S4. [Loop.] Go to step S1. ■

This is not, strictly speaking, an algorithm, since we didn't specify how the algorithm should choose the integers k and l that add up to m in step S3. However, running through a few instances with, say, $n = 6$, on paper, using whatever choices of split you like in every iteration of step S3, you'll find that the algorithm always terminates with $\text{BALANCE} = 15$. Trying a few different starting values of n might lead you to conjecture the following proposition, which we will prove by (strong) induction.

Proposition 17. *For a given input n , Algorithm S always outputs $\text{BALANCE} = n(n - 1)/2$, regardless of the choice of split at any given iteration of step S3.*

Proof. By induction on n . For the base case $n = 1$, note that we immediately output $\text{BALANCE} = 0$ in the very first step of the algorithm, and $0 = 1(1 - 1)/2$.

Now for the inductive step, let $n \geq 1$ and suppose that for $1 \leq k \leq n$, the payout for running algorithm S on input l is $l(l - 1)/2$. Suppose we have a pile of $n + 1$ stones. We shall divide it into piles of size k and size $n + 1 - k$. The total payout will be the pay for this division, namely $k(n + 1 - k)$, plus the pay for further subdividing the two piles. Thus by the induction hypothesis applied twice, the total payout will be

$$\begin{aligned} k(n + 1 - k) + \frac{k(k - 1)}{2} + \frac{(n + 1 - k)(n - k)}{2} \\ &= \frac{2nk + 2k - 2k^2 + k^2 - k + n^2 + n - 2nk - k + k^2}{2} \\ &= \frac{n^2 + n}{2} \\ &= \frac{(n + 1)((n + 1) - 1)}{2}, \end{aligned}$$

which is the expected formula for $n + 1$. ■

5. Functions

A *function* f from a set A to a set B is a subset $f \subseteq A \times B$ such that for every $a \in A$, there is exactly one $b \in B$ such that $(a, b) \in f$. (If there is no $b \in B$, or more than one, we say that f is not well-defined. If $(a, b) \in f$, we write $f(a) = b$. The set A is called the *domain* and the set B is the *codomain*. The notation $f : A \rightarrow B$ is a way of concisely writing “ f is a function with domain A and codomain B .”

Here are some examples and non-examples.

- i) The function $f : \mathbf{R} \rightarrow \mathbf{R}$ given by $f(x) = x^2$ is

$$f = \{(a, b) \in \mathbf{R}^2 : b = a^2\},$$

when written in set-builder notation.

- ii) On the other hand, the set

$$g = \{(a, b) \in \mathbf{R}^2 : a = b^2\}$$

is not a function, since $(1, 1)$ and $(1, -1)$ are both in g . Furthermore, there is no element of g with -1 as its first coordinate.

- iii) If $X = \{1, 2, 5\}$ and $Y = \{0, 1, 2, 3, 4, 5\}$, then

$$\{(1, 0), (2, 4), (5, 5)\}$$

is a function from X to Y .

- iv) The set

$$h = \{(x, y) \in \mathbf{R}^2 : y = 1/x\}$$

is not a function, since there is no y such that $(0, y) \in h$. However if we amend the domain and consider

$$h = \{(x, y) \in (\mathbf{R} \setminus \{0\}) \times \mathbf{R} : y = 1/x\},$$

then in fact, $h : \mathbf{R} \setminus \{0\} \rightarrow \mathbf{R}$ is a function.

24.IX **Injective and surjective functions.** The *range* or *image* of a function $f : A \rightarrow B$ is

$$f(A) = \{b \in B : \text{there exists } a \in A \text{ such that } b = f(a)\}.$$

These are all the values f actually outputs. For instance, if we let $f : \mathbf{Z} \rightarrow \mathbf{N}$ be given by $f(n) = n^2$, then

$$\begin{aligned} f(\mathbf{Z}) &= \{\dots, f(-2), f(-1), f(0), f(1), f(2), \dots\} \\ &= \{\dots, 9, 4, 1, 0, 1, 4, 9, \dots\} \\ &= \{0, 1, 4, 9, 16, 25, \dots\}. \end{aligned}$$

A function $f : A \rightarrow B$ is called *surjective* or *onto* if $f(A) = B$, that is, for every $b \in B$ there exists some $a \in A$ such that $f(a) = b$. An example of a surjective function is $f : \mathbf{Q} \rightarrow \mathbf{Q}$ sending $x \mapsto x/2$. This is because for any $q \in \mathbf{Q}$, we can set $r = 2q$, and

$$f(r) = \frac{r}{2} = \frac{2q}{2} = q.$$

A function $f : A \rightarrow B$ is called *injective* or *one-to-one* if for all $a_1, a_2 \in A$ with $a_1 \neq a_2$, we also have $f(a_1) \neq f(a_2)$. Equivalently, f is injective if $f(a_1) = f(a_2)$ implies that $a_1 = a_2$ for all $a_1, a_2 \in A$. For instance, the function $f : \mathbf{Z} \rightarrow \mathbf{N}$ that sends $n \mapsto n^2$ is not injective since $f(-1) = f(1)$, but $-1 \neq 1$. On the other hand, if we modify the domain, considering $f : \mathbf{N} \rightarrow \mathbf{N}$ sending $n \mapsto n^2$, then now f is injective, since if $f(m) = f(n)$, then $m^2 = n^2$, and there is only one positive integer that squares to any given integer, so $m = n$.

Hence we see that any function can be transformed into an injective one, in principle, by shrinking its domain (though this new function might no longer have the properties you liked in the original one), and any function $f : A \rightarrow B$ can be made surjective by changing its codomain to its range, i.e., letting $B = f(A)$.

The pigeonhole principle. We now take a brief pause to introduce one of the most fundamental laws in discrete mathematics, called the pigeonhole principle. We begin with the following intuitive theorem.

Theorem 18. *Let a_1, a_2, \dots, a_n be a finite sequence (repeats allowed) of real numbers. Let*

$$a = \frac{1}{n} \sum_{i=1}^n a_i$$

be the average value of the sequence and let m be the maximum value the sequence attains. Then $m \geq a$.

Proof. We have

$$a = \frac{1}{n} \sum_{i=1}^n a_i \leq \frac{1}{n} \sum_{i=1}^n m = \frac{mn}{n} = m. \quad \blacksquare$$

This theorem can be summed up in one sentence: *The maximum is at least the average.* Don't underestimate this theorem even though its proof was a one-liner! It is often used to prove highly nontrivial results. (As an exercise, prove the similar statement: *The minimum is at most the average.*) From here we are now equipped to prove (a general version of) the pigeonhole principle.

Theorem 19. *Let A and B be finite sets with $|A| = m$ and $|B| = n$. For every function $f : A \rightarrow B$, then there is some $b \in B$ such that there are at least $\lceil m/n \rceil$ elements $a \in A$ with $f(a) = b$.*

Proof. Enumerate $B = \{b_1, b_2, \dots, b_n\}$. For $1 \leq i \leq n$, let r_i be the number of $a \in A$ such that $f(a) = b_i$. This is a sequence that adds up to m , since every a in A maps to exactly one element of B . So the sequence has average m/n , and by the previous theorem, there must be some j such that $r_j \geq m/n$. But the r_i are all actually integers (since cardinalities of finite sets are integers), meaning that $r_j \geq \lceil m/n \rceil$. Letting $b = b_j$ completes the proof. \blacksquare

The reason this is called the pigeonhole principle is because of the following special case.

Corollary 20 (*Pigeonhole principle*). Let $n \geq 2$. If n pigeons nest in $n - 1$ holes, there is at least one hole that contains at least two pigeons.

Proof. Let A be the set of pigeons and B the set of pigeonholes. Let f be any function sending the set of pigeons to the set of holes. By the previous theorem, there is some hole with at least $\lceil n/(n-1) \rceil = 2$ pigeons in it. (This is because $1 < n/(n-1) < 2$ for all integers $n \geq 2$.) ■

Bijections. A function f is called *bijective* (or a *bijection*, or a *one-to-one correspondence*) if it is injective and surjective. Bijections are important because of the following proposition.

Proposition 21. Let A and B be finite sets. Then

- i) there exists a bijection $f : A \rightarrow B$ if and only if $|A| = |B|$; and
- ii) if $|A| = |B|$ and $f : A \rightarrow B$ then f is injective if and only if f is surjective.

Proof. Suppose $|A| = |B| = n$. Then choose an ordering a_1, \dots, a_n of A and an ordering b_1, \dots, b_n of B . Let $f(a_i) = b_i$ for all $1 \leq i \leq n$. By construction, this is a bijection, proving one direction of (i).

On the other hand, suppose $|A| \neq |B|$ (so we prove this direction by contraposition). If $A < B$, then f cannot be surjective, since the image of f has size at most $|A| < |B|$ (at least one element of B must be missed). If $A > B$, then $|A|/|B| > 1$, so by Theorem 19, there is some element $b \in B$ such that the number of $a \in A$ mapping to b is at least $\lceil |A|/|B| \rceil \geq 2$. This means that f is not injective. We have proved part (i).

To prove part (ii), let $|A| = |B|$ and let $f : A \rightarrow B$. First we assume that f is injective. We enumerate $A = \{a_1, a_2, \dots, a_n\}$. Then

$$f(A) = \{f(a_1), f(a_2), \dots, f(a_n)\} \subseteq B.$$

All of the $f(a_i)$ are distinct, since if $f(a_i) = f(a_j)$, then $a_i = a_j$. So $|f(A)| = |A| = n$, and $f(A)$ is a size n subset of B , which also has size n . Hence $f(A) = B$; that is, f is surjective.

Lastly, suppose f is not injective (again we are using contraposition). So there are a_i and a_j such that $a_i \neq a_j$ but $f(a_i) = f(a_j)$. So

$$|f(A)| = |\{f(a_1), f(a_2), \dots, f(a_n)\}| < n = |B|,$$

so $f(A) \neq B$ and f is not surjective. ■

Item (i) of the previous proposition should be entirely intuitive, especially if we use the alternative nomenclature “one-to-one correspondence” instead of “bijection.” (In fact, we already implicitly used (i) in these notes, in the proof of Proposition 5.) Item (ii) is perhaps not as immediate, but should become clear if you work it out with a picture.

Bijections. A function $f : A \rightarrow B$ is called *invertible* if there exists $g : B \rightarrow A$ such that

- i) for all $b \in B$, $f(g(b)) = b$; and
- ii) for all $a \in A$, $g(f(a)) = a$.

If g exists, it can be shown that g must be unique, so we write $g = f^{-1}$ and speak of *the* inverse of f .

Proposition 22. *Let $f : A \rightarrow B$ be a function. Then f is invertible if and only if f is bijective.*

Proof. First we assume that f is invertible. So there exists an inverse g of f . For each $b \in B$, setting $a = g(b)$ we have

$$f(a) = f(g(b)) = b.$$

This proves that f is surjective. To show that f is injective, suppose that $f(a_1) = f(a_2)$. By applying g on both sides, we have $g(f(a_1)) = g(f(a_2))$, whence $a_1 = a_2$, by definition of g .

Now assume that f is bijective. We need to define $g : B \rightarrow A$. Well, given any $b \in B$, there is some a such that $f(a) = b$, from surjectivity of f , and this a is unique, since f is injective. So set $g(b) = a$ (and repeat this process for every $b \in B$). We have $f(g(b)) = f(a) = b$, and for every $a \in A$, by definition of g the element $g(f(a))$ is the unique element in A that gets brought to $f(a)$ by f , has to be a itself. ■

Sometimes to prove that two sets have the same cardinality, it is easier to prove that there exists a bijection (as we already saw in the example of Proposition 5), and sometimes to prove that a function is a bijection, it is easier to show that it has an inverse, rather than messing around with the definitions of injective and surjective. Here's an example.

Proposition 23. *Let X be a finite nonempty set. Let E be the set of all subsets of X with even cardinality, and let D be the set of all subsets of X with odd cardinality. Then $|E| = |D|$.*

Proof. We shall construct a function $f : E \rightarrow D$. Fix one specific $x \in X$; we can do this because $X \neq \emptyset$. Now, for all $A \in E$, let

$$f(A) = \begin{cases} A \setminus \{x\}, & \text{if } x \in A; \\ A \cup \{x\}, & \text{if } x \notin A. \end{cases}$$

Note that since $|A|$ is even for all $A \in E$, the cardinality of $f(A)$ is odd (in the first case it is $|A| - 1$ and in the second case it is $|A| + 1$). This shows that f is indeed a function with codomain D . Now to prove $|E| = |D|$ we will show that f is bijective, which we do by showing that f is an inverse (as an exercise, you might instead try to prove bijectivity from the definitions of injective and surjective).

We define $g : D \rightarrow E$ by

$$g(A) = \begin{cases} A \setminus \{x\}, & \text{if } x \in A; \\ A \cup \{x\}, & \text{if } x \notin A. \end{cases}$$

As before, $|g(A)|$ is even, since A is assumed to be a member of D now. Now for any $A \in E$,

$$\begin{aligned} g(f(A)) &= \begin{cases} g(A \setminus \{x\}), & \text{if } x \in A; \\ g(A \cup \{x\}), & \text{if } x \notin A \end{cases} \\ &= \begin{cases} (A \setminus \{x\}) \cup \{x\}, & \text{if } x \in A; \\ (A \cup \{x\}) \setminus \{x\}, & \text{if } x \notin A \end{cases} \\ &= \begin{cases} A, & \text{if } x \in A; \\ A, & \text{if } x \notin A \end{cases} \\ &= A. \end{aligned}$$

The proof that $f(g(A)) = A$ is similar. Thus g is the inverse of f . ■

6. Cardinality

26.IX Earlier, we defined the cardinality of a set to be the number of elements it contains. What, then, is the cardinality of \mathbf{N} ? How about \mathbf{R} ? You might say ∞ , but this is not a number (at least, it's not an element of \mathbf{N} , the way all cardinalities of finite sets are). So perhaps we should amend our question to the following: *When do infinite sets have the same size?* Our experience with functions leads us to the answer: *When there exists a bijection between them.* We shall say that A and B are *equipotent* (or *equinumerous*, or *have the same cardinality*) if there exists a bijection between A and B . In this case we write $|A| = |B|$.

As an example, the sets \mathbf{N} and $\mathbf{N} \setminus \{0\}$ are equipotent, since f given by $n \mapsto n + 1$ is a bijection $\mathbf{N} \rightarrow \mathbf{N} \setminus \{0\}$. (Check that $f^{-1}(m) = m - 1$ is its inverse.)

It is even possible to remove an infinite number of elements from \mathbf{N} and end up with something still equipotent with \mathbf{N} . To see this, let E be the set of nonnegative even integers, and consider the function $f : \mathbf{N} \rightarrow E$ sending $n \mapsto 2n$. This is injective because if $2m = 2n$, then dividing out by 2 on both sides yields $m = n$. It is surjective because if $n \in E$, then $n = 2k$ for some $k \in \mathbf{N}$, by definition, and $f(k) = 2k = n$.

So we can find subsets of \mathbf{N} equipotent with it. It turns out we can also find supersets of \mathbf{N} with the same property.

Theorem 24. We have $|\mathbf{N}| = |\mathbf{Z}|$.

Proof. We define $f : \mathbf{N} \rightarrow \mathbf{Z}$ by

$$f(n) = \begin{cases} \frac{n}{2}, & \text{if } n \text{ is even;} \\ -\frac{n+1}{2}, & \text{if } n \text{ is odd.} \end{cases}$$

We shall show that f is a bijection.

Note first that if n is even, then $f(n) \geq 0$, and if n is odd, then $f(n) < 0$. So if $f(m) = f(n)$ for some $m, n \in \mathbf{N}$, then $f(m)$ and $f(n)$ must either both be

negative, or both be nonnegative. Either way, m and n are either both even or they are both odd. If m and n are both even, then from $f(m) = f(n)$ we derive

$$\frac{m}{2} = \frac{n}{2},$$

whence multiplying by 2 on both sides we see that $m = n$. If m and n are both odd, then

$$-\frac{m+1}{2} = -\frac{n+1}{2},$$

so, multiplying by -2 and subtracting 1 from both sides we have $m = n$ in this case as well.

Now we show that f is surjective. Let $k \in \mathbf{Z}$. If $k \geq 0$, then consider $n = 2k$. We have

$$f(n) = f(2k) = \frac{2k}{2} = k.$$

If $k < 0$, then consider $n = -2k - 1$. (Check that this is an element of \mathbf{N} .) Then

$$f(n) = f(-2k - 1) = -\frac{-2k - 1 + 1}{2} = -\frac{-2k}{2} = k.$$

This shows that f is surjective, so f is in fact bijective and we conclude that $|\mathbf{N}| = |\mathbf{Z}|$. ■

We say that a set A is *countably infinite* if there exists a bijection $f : \mathbf{N} \rightarrow A$, that is, if $|\mathbf{N}| = |A|$. A set is said to be *countable* if it is either finite or countably infinite. Otherwise it is called *uncountable*. The previous theorem shows that \mathbf{Z} is countably infinite.

Sometimes it is difficult to come up with a bijection directly. Instead, we would like to find an injection from A to B (which, in some sense, shows that $|A| \leq |B|$), and then an injection from B to A . This is made possible by the following useful theorem, named for Ernst Schröder and Felix Bernstein, who independently proved it in 1898. The proof is a bit difficult, so it's technically outside the scope of the course. For fun, you might try to do it as an exercise.

Theorem 25. (*Schröder–Bernstein theorem*). *If there exists an injective function $f : A \rightarrow B$ and another injective function $g : B \rightarrow A$, then there is a bijection $h : A \rightarrow B$.*

**Proof.* We present the proof as a (difficult) exercise. Here is the roadmap. Call $b \in B$ *unattached* if there is no $a \in A$ such that $f(a) = b$. Let $h : B \rightarrow B$ be given by $h(b) = f(g(b))$. Given $b, b' \in B$, we shall say that b is a *peer* of b' if either $b = b'$ or there exists some $n \in \mathbf{N}$ such that

$$b = \underbrace{h(h(\cdots(h(b'))\cdots))}_{n \text{ times.}}$$

Say that $b \in B$ is a *PAE* if it is the peer of an unattached element. (So unattached elements are automatically PAEs, by setting $n = 0$.)

- a) Show that if $a \in A$ is such that $f(a)$ is a PAE, then there is a unique element $b^* \in B$ such that $g(b^*) = a$, and that this element is a PAE.

By part (a), if $f(a)$ is a PAE, it makes sense to speak of $g^{-1}(a)$. It is the element $b^* \in B$ such that $g(b^*) = a$. From here we define

$$r(a) = \begin{cases} g^{-1}(a), & \text{if } f(a) \text{ is a PAE;} \\ f(a), & \text{otherwise.} \end{cases}$$

- b) Show that if $b \in B$ is a PAE then so is $f(g(b))$.
 c) Show that r is surjective. [*Hint*: Do a proof by cases. Every $b \in B$ is either a PAE or it is not a PAE.]
 d) We already proved in part (a) that if $f(a)$ is a PAE, then so is $r(a)$. Prove the converse of this statement.
 e) Show that r is injective, and therefore bijective. [*Hint*: Assume $r(a_1) = r(a_2)$, and do a proof by cases again. Part (d) will be useful here.] ■

Once again, this proof is outside the scope of the course. Don't lose sleep over it if you can't do it. Feel free to ask questions at office hours if you get stuck.

Using the Schröder–Bernstein theorem, we now show that the Cartesian product of two countable sets is also countable. In the proof, we shall employ the fact that if a set A is countable, then there exists an enumeration

$$A = \{a_0, a_1, a_2, \dots\}.$$

(If f is the bijection given by the definition, we can set $a_0 = f(0)$, $a_1 = f(1)$, and so on.)

In the following proof, we'll also use the Fundamental Theorem of Arithmetic, which we state now, and prove later in the course (in the number theory section). It says that any integer can be factored as a product of primes, a theorem you should have learned in grade school.

Theorem 26 (*Fundamental Theorem of Arithmetic*). Every positive integer $n \geq 2$ can be factored into a product

$$p_1^{v_1} p_2^{v_2} \cdots p_m^{v_m},$$

where $m \geq 0$ is an integer, p_1, p_2, \dots, p_m are distinct primes, and v_1, v_2, \dots, v_m are positive integers. This factorisation is unique up to the order of the primes. ■

Theorem 27. If A and B are countably infinite sets, then $A \times B$ is also countably infinite.

Proof. Enumerate $A = \{a_0, a_1, \dots\}$ and $B = \{b_0, b_1, \dots\}$. Now given $(a_i, b_j) \in A \times B$ for some $i, j \in \mathbb{N}$, we can let

$$f(a_i, b_j) = 2^i 3^j.$$

We show that this defines an injective function. Suppose that $f(a_i, b_j) = f(a_{i'}, b_{j'})$. Then $2^i 3^j = 2^{i'} 3^{j'}$, and by the uniqueness of prime factorisations in the Fundamental Theorem of Arithmetic, we see that $i = i'$ and $j = j'$.

Now we produce an injection $g : \mathbf{N} \rightarrow A \times B$ by simply setting

$$g(n) = (a_n, b_0).$$

It is clear this is injective, since if $g(m) = g(n)$, then $(a_m, b_0) = (a_n, b_0)$, and that implies that $m = n$. ■

Corollary 28. *We have $|\mathbf{Z} \times \mathbf{Z}| = |\mathbf{N}|$.*

Proof. We proved earlier that \mathbf{Z} is countably infinite, so we may apply the previous theorem with $A = B = \mathbf{Z}$. ■

This corollary concerning $\mathbf{Z} \times \mathbf{Z}$ allows us to prove that \mathbf{Q} is countable as well.

Theorem 29. *The set \mathbf{Q} of rational numbers is countable.*

Proof. We define $f : \mathbf{Q} \rightarrow \mathbf{Z} \times \mathbf{Z}$ as follows. For an element $q \in \mathbf{Q}$, we let $q = a/b$ the fraction q written in lowest terms (where $a \in \mathbf{Z}$ and $b \in \mathbf{N} \setminus \{0\}$). Then we set $f(q) = (a, b)$. This is an injective function because if $f(q) = f(q')$, then writing $q = a/b$ and $q' = a'/b'$ in lowest terms, we have $(a, b) = (a', b')$, so $q = a/b = a'/b' = q'$.

To define an injection $g : \mathbf{Z} \times \mathbf{Z} \rightarrow \mathbf{Q}$, we recycle the injection we had from the proof of Theorem 27. We already know that \mathbf{Z} is countable, so fix an enumeration $\mathbf{Z} = \{a_0, a_1, a_2, \dots\}$. Then let $g(a_i, a_j) = 2^i 3^j$. The range of g is a subset of \mathbf{N} , so *a fortiori* it is a subset of \mathbf{Q} , and we already showed before that it is an injection.

We have shown that $|\mathbf{Q}| = |\mathbf{Z} \times \mathbf{Z}|$, which in turn shows that $|\mathbf{Q}| = |\mathbf{N}|$, after applying Corollary 28. ■

So far, we have just given lots of examples of countably infinite sets. Finally, we give an example of a set that is not countably infinite.

Theorem 30. *The set A of all infinite binary strings is uncountable.*

Proof. Certainly A is not finite, so we need to show that $|\mathbf{N}| \neq A$.

Let $f : \mathbf{N} \rightarrow A$. We shall show that f is not surjective. For all $m, n \in \mathbf{N}$, let $a_{m,n}$ be the n th bit of $f(m)$. Consider the infinite binary string

$$s = (1 - a_{0,0}, 1 - a_{1,1}, 1 - a_{2,2}, 1 - a_{3,3}, \dots).$$

This string s cannot equal $f(m)$ for any $m \in \mathbf{N}$, since given an arbitrary m , the m th bit of $f(m)$ is $a_{m,m}$, whereas the m th bit of s is $1 - a_{m,m}$. Hence there is some $s \in A$ such that $f(m) \neq s$ for all $m \in \mathbf{N}$. So f is not surjective.

Since there can be no surjection $f : \mathbf{N} \rightarrow A$, *a fortiori* we cannot have a bijection $\mathbf{N} \rightarrow A$. We conclude that $|A| \neq |\mathbf{N}|$. ■

The proof above was published by Georg Cantor in 1891, and hence is known as *Cantor's diagonal argument*. The technique has since been used to prove many other things.

The set of all infinite binary strings is in bijection with the set of all real numbers in the interval $[0, 1]$. If we're being extra pedantic, we need to forbid an infinite trailing string of all 1s, since, e.g., $0.0111\dots = 0.1000\dots$ after carrying the 1s (akin to the fact that, e.g., $0.0999\dots = 0.1$). After dealing with this detail, one will have shown that $|[0, 1]| \neq |\mathbf{N}|$.

7. Relations

01.X A *relation* on a set X is a subset $R \subseteq X \times X$. If $(a, b) \in R$ we write aRb and say “ a is related to b .” Here are some examples.

- i) The set $L = \{(a, b) \in \mathbf{R} \times \mathbf{R} : a < b\}$ is a relation. For instance, we can write $(2, \pi) \in L$, or $2L\pi$, or $2 < \pi$ (which is the more common notation for this relation).
- ii) The set $E = \{(a, b) \in \mathbf{R} \times \mathbf{R} : a = b\}$ is a relation (it is the “equals” relation).
- iii) On \mathbf{Z} , the set

$$R = \{(-1, 4), (8, -3), (0, 0), (0, 1)\}$$

is a relation.

- iv) For any set A , a function $f : A \rightarrow A$ is by definition a subset of $A \times A$, and hence is an example of a relation.
- v) Let H be the set of all humans, define $M \subseteq H \times H$ by setting (h_1, h_2) if and only if h_1 is married to h_2 . For instance,

$$(\text{Michelle Obama}, \text{Barack Obama}) \in M.$$

On the other hand, it is unfortunately the case that

$$(\text{Marcel Goh}, \text{Taylor Swift}) \notin M.$$

Properties of relations. Let $R \subseteq A \times A$ be a relation. R is called

- i) *reflexive* if for all $a \in A$, aRa ;
- ii) *symmetric* if for all $a, b \in A$ with aRb , we also have bRa ; and
- iii) *transitive* if for all $a, b, c \in A$ with aRb and bRc , we also have aRc .

Going back to the examples we had above, the relation L is transitive but neither reflexive nor symmetric, the relation E satisfies all three properties, the relation R satisfies none of them, and the relation M is symmetric but neither reflexive nor transitive. (Convince yourself of all of these facts.)

As a more involved example, let X be any set with $|X| \geq 2$. On 2^X , define a relation R by

$$(A, B) \in R \Leftrightarrow A \cap B \neq \emptyset$$

for all $A, B \subseteq X$. Which of the three properties above does R satisfy?

Is it reflexive? Well, is it true that for all $A \subseteq X$, $A \cap A \neq \emptyset$? The answer is no, since we have $\emptyset \cap \emptyset = \emptyset$, so we have $(\emptyset, \emptyset) \notin R$.

Is R symmetric? Well, if $A \cap B \neq \emptyset$, then since \cap is commutative, we have $B \cap A \neq \emptyset$ as well, so $(A, B) \in R$ implies that $(B, A) \in R$. In other words, yes, R is symmetric.

Is R transitive? If $A \cap B \neq \emptyset$ and $B \cap C \neq \emptyset$, does that necessarily mean that $A \cap C \neq \emptyset$? The answer is no. Here's the proof. Since $|X| \geq 2$ we can find $x, y \in X$ with $x \neq y$. Let $A = \{x\}$, $B = \{x, y\}$, and $C = \{y\}$. Then $A \cap B = \{x\} \neq \emptyset$, $B \cap C = \{y\} \neq \emptyset$, but alas $A \cap C = \emptyset$.

Equivalence relations and classes. If R is reflexive, symmetric, and transitive, then we say that R is an *equivalence relation*. Here are two examples (as an exercise, prove that both of these are equivalence relations).

- i) Let F be the set of all formulas in propositional logic with the variables p , q , and r , and operations \neg , \wedge , and \vee . The relation \equiv defined by $f_1 \equiv f_2$ if and only if f_1 and f_2 have the same truth table is an equivalence relation.
- ii) On the set \mathbf{R}^2 , define the relation $R \subseteq \mathbf{R}^2 \times \mathbf{R}^2$ by letting $(x, y)R(w, z)$ if and only if $\sqrt{x^2 + y^2} = \sqrt{w^2 + z^2}$. This is also an equivalence relation

If R is an equivalence relation, often we shall write $x \sim y$ to mean xRy . Sometimes we might even just say that \sim is the equivalence relation.

Let $R \subseteq A \times A$ be an equivalence relation. Define the *equivalence class* of $a \in A$ to be the set

$$[a] = \{b \in A : a \sim b\}.$$

This is the set of all $b \in A$ that are related to a .

Take for instance the example F above of all formulas under the relation \equiv . The equivalence class of the formula $p \vee \neg p$ is the set of all formulas whose truth table contains only 1s, that is, the set of all tautologies. And for the example R above, the equivalence class of the point $(1, 3)$ is the set

$$[(1, 3)] = \{(x, y) \in \mathbf{R}^2 : \sqrt{10} = \sqrt{x^2 + y^2}\};$$

that is, the circle of radius $\sqrt{10}$ centred about the origin.

We have the following proposition concerning equivalence relations.

Proposition 31. *Let R be an equivalence relation on A . Then*

- i) *for all $x \in A$, $x \in [x]$;*
- ii) *for all $x, y \in A$, $x \sim y$ if and only if $[x] = [y]$; and*
- iii) *for all $x, y \in A$, $x \not\sim y$ if and only if $[x] \cap [y] = \emptyset$.*

Proof. Let $x \in A$. Since R is reflexive, $x \sim x$, so we have $x \in [x]$. This proves part (i).

For part (ii), first we prove the “only if” direction. Let $x, y \in A$ be such that $x \sim y$. To prove that $[x] \subseteq [y]$, we let $z \in [x]$; so $x \sim z$. By symmetry, we have $z \sim x$, and this combined with $x \sim y$ allow us to deduce that $z \sim y$, by transitivity. Then by symmetry again, we have $y \sim z$, so $z \in [y]$. *Mutatis mutandis*, i.e., by swapping the roles of x and y , we also have $[y] \subseteq [x]$. Hence $[x] = [y]$.

Now for the “if” direction of (ii). Let $x, y \in A$ be such that $[x] = [y]$. By (i), we have $x \in [x]$, but then this means that $x \in [y]$. By definition of $[y]$, this means that $y \sim x$, so $x \sim y$ by symmetry.

On to part (iii). We prove both implications by contraposition (that is, we negate both sides of the statement). Let $x, y \in A$ be such that $[x] \cap [y] \neq \emptyset$. This means there is some $z \in A$ such that $z \in [x]$ and $z \in [y]$; so $x \sim z$ and $y \sim z$. By symmetry, $z \sim y$, so by transitivity, we have $x \sim y$.

On the other hand, suppose that $x, y \in A$ satisfy $x \sim y$. By (ii), we have $[x] = [y]$, so $[x] \cap [y] = [x] \neq \emptyset$, where we know that $[x] \neq \emptyset$ because it contains at least x (again, using (i)). ■

Let A be a set. A *partition* of A is a set P of subsets of A (i.e., $P \subseteq 2^A$), such that

- i) for all $x \in A$ there exists $S \in P$ such that $x \in S$;
- ii) for all $S_1, S_2 \in P$ with $S_1 \neq S_2$, the intersection $S_1 \cap S_2$ is empty; and
- iii) $\emptyset \notin P$.

Here are some examples.

- i) Let E be the set of all even integers and F the set of all odd integers. Then $\{E, F\}$ is a partition of \mathbf{Z} .
- ii) The set $\{(-\infty, 0), \{0\}, (0, \infty)\}$ is a partition of \mathbf{R} .

03.X If \sim is an equivalence relation on a set A , then we can define the *quotient* of A by \sim as the set of all equivalence classes of A under \sim . We denote this set by

$$A/\sim = \{[x] : x \in A\}.$$

We use the previous proposition to prove that quotients of sets by equivalence relations are partitions.

Proposition 32. *Let A be a set and \sim an equivalence relation on A . Then A/\sim is a partition of A .*

Proof. By part (i) of the previous proposition, every $x \in A$ belongs to the equivalence class $[x]$. Then, by part (ii) of the previous proposition, we know that for any equivalence classes $[x]$ and $[y]$ such that $[x] \neq [y]$, we must have $x \not\sim y$, and by part (iii) of the previous proposition, we deduce that $[x] \cap [y] = \emptyset$. This shows that A/\sim satisfies the second part of the definition of partition. Lastly, we note

that $\emptyset \notin A/\sim$, since every element of A/\sim is equal to $[x]$ for some $x \in A$, and must thus contain at least the element x . ■

Let us now revisit the examples of equivalence relations from last class, and see what partitions they give rise to. In the example F of propositional formulas with variables p , q , and r , under the equivalence relation \equiv , the set of equivalence classes is the set of all possible truth tables on three variables. Each such truth table has 8 rows, so there are 2^8 equivalence classes. In other words, $|F/\equiv| = 2^8$.

How about the relation R defined on \mathbf{R}^2 where

$$(x, y) \sim (w, z) \quad \text{if and only if} \quad \sqrt{x^2 + y^2} = \sqrt{w^2 + z^2}?$$

Well, each $[(x, y)]$ is the circle of radius $\sqrt{x^2 + y^2}$ centred around the origin, so \mathbf{R}^2/\sim is the set of all circles in the plane centred at $(0, 0)$.

II. NUMBER THEORY

*Die Mathematik ist die Königin der Wissenschaften
und die Zahlentheorie ist die Königin der Mathematik.*

— Attributed to C. F. GAUSS (1777–1855)

8. Division

Let $a, b \in \mathbf{Z}$. We say that a *divides* b (or b is a *multiple* of a , or b is *divisible* by a , or a is a *factor* of b) if there exists $n \in \mathbf{Z}$ such that $b = na$. In this case we write $a \mid b$. For example, $2 \mid 10$, since $10 = 5 \cdot 2$, but 3 does not divide 10, since there does not exist $n \in \mathbf{Z}$ such that $10 = 3n$. This defines a relation on \mathbf{Z} .

Note that for all $n \in \mathbf{Z}$, $n \mid 0$, since $0 = 0 \cdot n$. We also have $1 \mid n$ for all $n \in \mathbf{Z}$; since $n = n \cdot 1$. It is true that $0 \mid 0$, since we have, say, $0 = 1 \cdot 0$, but for all nonzero $n \in \mathbf{N}$, 0 does not divide n , since all multiples of 0 equal 0 (and thus cannot equal n). Further properties of the “divides” relation are given by the next proposition.

Proposition 33. *For all $a, b, c, d \in \mathbf{Z}$,*

- i) *if $a \mid b$, then $a \mid bc$;*
- ii) *if $a \mid b$ and $a \mid c$, then $a \mid (b + c)$;*
- iii) *if $a \mid b$ and $b \mid c$, then $a \mid c$;*
- iv) *if $a \mid b$ and $b \neq 0$, then $|a| \leq |b|$; and*
- v) *if $a \mid b$ and $b \mid a$, then $|a| = |b|$.*

Proof. We leave parts (i) and (ii) as exercises to the reader.

For part (iii), if $a \mid b$ and $b \mid c$, then there are integers k and l such that $b = ka$ and $c = lb$. Then $c = kla$, so $a \mid c$ (since kl is also an integer).

For part (iv), suppose that $a \mid b$ and $b \neq 0$. Then $b = ka$ for some $k \in \mathbf{Z}$, and $k \neq 0$ since $b \neq 0$. This means that $|b| = |k| \cdot |a|$, but $|k| \geq 1$, so $|b| \geq |a|$.

For part (v), assume that $a \mid b$ and $b \mid a$. If $a \neq 0$ and $b \neq 0$, then we may apply part (iv) twice to get $|a| \leq |b|$ and $|b| \leq |a|$, which together imply $|a| = |b|$. If $b = 0$, then $0 \mid a$ so $a = 0$ as well, and $|a| = |b|$. Likewise, if $a = 0$, then $0 \mid b$, so $b = 0$ and in this case as well, $|a| = |b|$. ■

In grade school, you learned how to divide an integer by another one, obtaining a quotient and a remainder. We state this as a theorem. Its proof (which we shall consider outside the scope of the course, but which we include as optional reading for those interested) relies on the well-ordering principle (which we’ve used already in these notes) and the *Archimedean property* of \mathbf{R} , which states that for every $x \in \mathbf{R}$ there exists $n \in \mathbf{N}$ such that $n > x$. (One can learn the proof of the Archimedean property from an introductory course in analysis, e.g., MATH 242/254.)

Theorem 34 (*Division algorithm*). *Let $a, b \in \mathbf{Z}$ with $b > 0$. Then there exist unique integers q and r such that $a = bq + r$ and $0 \leq r < b$.*

**Proof.* First we show that such integers q and r exist. If $b \mid a$ then $a = kb$ for some $k \in \mathbf{Z}$, and we can set $q = k$ and $r = 0$.

If a is not divisible by b , then consider the numbers

$$\dots, a - 3b, a - 2b, a - b, a, a + b, a + 2b, a + 3b, \dots$$

Let S be the set of these integers that are positive. Symbolically, we have

$$S = \{a - kb : k \in \mathbf{Z} \text{ and } a - kb \geq 0\} \quad .$$

By the Archimedean property, there is some $n \in \mathbf{N}$ such that $n > -a$, which implies that $nb \geq n > -a$ (here we use the fact that if $b > 0$ and b is an integer, then $b \geq 1$). From this we derive $a > -nb$, and hence $a + nb = a - (-n)b > 0$. This shows that S is nonempty.

Since S is a nonempty subset of \mathbf{N} , by the well-ordering principle it has a least element, call it r . By definition of S , there must be some integer q such that $r = a - qb$, so $a = bq + r$. We now claim that $0 < r < b$.

We know that $r > 0$, since all elements of S are positive by definition. Suppose, for a contradiction, that $r \geq b$. Then $a - bq = r \geq b$, and so

$$0 \leq r - b = a - qb - b = a - (q + 1)b.$$

Since $q + 1$ is an integer, by definition of S , either $r - b$ is an element of S , or $r - b = 0$. Since r was defined to be the minimal element of S , it cannot be the case that $r - b$ is in S . So $r - b = 0$. But this means that $0 = a - (q + 1)b$; that is, $a = (q + 1)b$, contradicting our assumption that b does not divide a . The contradiction allows us to conclude that $0 < r < b$ (in the case that b does not divide a). In general, we have shown that $0 \leq r < b$.

Lastly, we need to prove that q and r are uniquely determined by the integers a and b . Suppose that

$$a = bq_1 + r_1 \quad \text{and} \quad 0 \leq r_1 < b$$

and

$$a = bq_2 + r_2 \quad \text{and} \quad 0 \leq r_2 < b,$$

for some integers q_1, q_2, r_1 , and r_2 . We shall show that $q_1 = q_2$ and $r_1 = r_2$. Well, suppose that $r_1 \neq r_2$, for a contradiction. Without loss of generality we can assume that $r_1 < r_2$. Then subtracting the two equations, we obtain

$$0 = a - a = (bq_1 + r_1) - (bq_2 + r_2) = b(q_1 - q_2) + (r_1 - r_2).$$

This means that

$$r_2 - r_1 = b(q_1 - q_2),$$

so we find that $b \mid (r_2 - r_1)$. By part (iv) of the previous proposition, we obtain $|b| \leq |r_2 - r_1|$, and we can simply write $b \leq r_2 - r_1$, since both of these quantities are positive. But this is a contradiction, since

$$0 \leq r_1 < r_2 < b,$$

yields $r_2 - r_1 < b$. The contradiction shows that that $r_2 = r_1$. Substituting this into the relation $r_2 - r_1 = b(q_1 - q_2)$, we get $0 = b(q_1 - q_2)$ and conclude that $q_1 - q_2 = 0$, since $b > 0$. ■

Let a and b be integers, not both zero. Their *greatest common divisor*, $\gcd(a, b)$ is defined to be the greatest positive integer d such that $d \mid a$ and $d \mid b$. Note that $\gcd(0, 0)$ is not defined, since all positive integers d satisfy $d \mid 0$. On the other hand $\gcd(x, 0)$ is simply $|x|$, and $\gcd(x, 1) = 1$. Lastly, we don't need to worry about negative signs when computing greatest common divisors; i.e., $\gcd(\pm x, \pm y) = \gcd(|x|, |y|)$.

Euclid's algorithm. Now we ask ourselves, How do we compute greatest common divisors in general? The answer lies in one of the oldest algorithms known to humankind. It appears in Euclid's *Elements*, written around 300 B.C.

Algorithm E (*Euclid's algorithm*). Given two nonnegative integers a and b , not both zero, this algorithm outputs $\gcd(a, b)$.

E1. If $b = 0$, then output a and terminate the algorithm.

E2. Since $b \neq 0$, by the division algorithm we may write $a = qb + r$, where $0 \leq r < b$. Set $a \leftarrow b$, $b \leftarrow r$, and return to step E1. ■

The algorithm will eventually terminate, since the stopping criterion is that b be equal to 0, and in step E2 we replace b with a number that is strictly closer to 0. But will it terminate with the correct answer? Well, we know step E1 is correct, because of our earlier observation that $\gcd(a, 0) = a$ (whenever a is positive). On the other hand, it is not at all evident that step E2 will eventually output $\gcd(a, b)$, since we actually overwrite the values of a and b in the step! The following lemma clarifies the situation.

Lemma 35. Let $a, b, q, r \in \mathbf{Z}$ be integers such that $a = qb + r$. Then $\gcd(a, b) = \gcd(b, r)$.

Proof. We shall show that for any $d \in \mathbf{Z}$,

$$d \mid a \text{ and } d \mid b \quad \text{if and only if} \quad d \mid b \text{ and } d \mid r.$$

For the forward implication, suppose that $a = kd$ and $b = ld$ for some $k, l \in \mathbf{Z}$. Substituting this into the identity $a = qb + r$ yields $kd = ldq + r$, whence $r = d(k - lq)$, so we conclude that $d \mid r$. (This is because $k - lq \in \mathbf{Z}$.)

For the reverse implication, suppose that $b = ld$ and $r = md$ for some $k, l \in \mathbf{Z}$. Substituting this into $a = qb + r$, we have $a = ldq + md$, so $a = d(lq + m)$, which means that $d \mid a$.

We have proved that a and b have the same common factors as b and r , so they must have the same greatest common divisor. ■

Now that we are secure in the fact that Euclid's algorithm will indeed terminate with the correct output, let us now see it in action. Suppose we want to find $\gcd(30, 112)$. We perform successive divisions replacing the pair (a, b) with

a new pair (b, r) each time:

$$\begin{aligned}
 30 &= 0 \cdot 112 + 30 \\
 112 &= 3 \cdot 30 + 22 \\
 30 &= 1 \cdot 22 + 8 \\
 22 &= 2 \cdot 8 + 6 \\
 8 &= 1 \cdot 6 + 2 \\
 6 &= 3 \cdot 2 + 0
 \end{aligned} \tag{3}$$

We stop once the remainder r equals 0, and the answer is $\gcd(b, r) = \gcd(b, 0) = b$. (So in the example above, the final answer is 2.)

Bézout's identity. The following theorem allows us to express the greatest common divisor as a linear combination of the two integers in question.

Theorem 36 (*Bézout's identity*). *Let a and b be nonzero integers with greatest common divisor $\gcd(a, b)$. Then there exist integers s and t such that*

$$\gcd(a, b) = sa + tb.$$

Moreover, $\gcd(a, b)$ is the least positive integer that can be expressed as an integer linear combination of a and b .

o8.x *Proof.* Let

$$S = \{s'a + t'b : s', t' \in \mathbf{Z}, ax + by > 0\}.$$

This set is nonempty, since if a is negative then $(-1)a + 0b \in S$ and if a is positive, then $1a + 0b \in S$. Since S is a nonempty set of positive integers, it has a least element, by the well-ordering principle. Call this integer $d = sa + tb$ (for some specific choices of $s, t \in \mathbf{Z}$); the claim is that $d = \gcd(a, b)$.

By the division algorithm, we may write

$$a = dq + r$$

where q and r are integers with $0 \leq r < d$. But we can write

$$\begin{aligned}
 r &= a - qd \\
 &= a - q(sa + tb) \\
 &= (1 - qs)a + (qt)b,
 \end{aligned}$$

so $r \in S \cup \{0\}$. But $r < d$, so if $r \in S$, then d would not be the smallest element of S . So we must have $r = 0$. This implies that d is a divisor of a . Repeating this argument with b instead of a , we find that d divides b as well.

We have shown that d is a common divisor of a and b . It remains to show that it is the greatest one. That is, we must show that if $c \mid a$ and $c \mid b$, then $c \leq d$. Well, if $a = kc$ and $b = lc$, then the identity

$$d = sa + tb = skc + tlc = (sk + tl)c,$$

shows that d is a multiple of c as well. Since $d > 0$, this means that $c \leq d$. ■

This theorem is named for Étienne Bézout, who proved an analogous result (with polynomials instead of integers) in 1779, but the result above for integers has been known since at least the 1600s.

To actually find the integers s and t such that $\gcd(a, b) = sa + tb$, we first perform the Euclidean algorithm, keeping track of all our intermediate steps. Then we combine all the information from each step to work out what s and t are. For example, in the earlier example showing that $2 = \gcd(112, 30)$, we start with

$$2 = 8 - 1 \cdot 6,$$

which is (a rearrangement of) the fifth line of (3). Then the fourth line of (3) says that $6 = 22 - 2 \cdot 8$, so

$$2 = 8 - 1 \cdot (22 - 2 \cdot 8) = 8 - 22 + 2 \cdot 8 = 3 \cdot 8 - 22.$$

The third line of (3) tells us that $8 = 30 - 1 \cdot 22$, which gives

$$2 = 3 \cdot (30 - 1 \cdot 22) - 22 = 3 \cdot 30 - 3 \cdot 22 - 22 = 3 \cdot 30 - 4 \cdot 22.$$

We now have the number 30 appearing in the expression, we just need to get rid of the 22 and replace it with 112. To do this we use the second line of (3), which says that $22 = 112 - 3 \cdot 30$. We end up with

$$2 = 3 \cdot 30 - 4(112 - 3 \cdot 30) = 3 \cdot 30 - 4 \cdot 112 + 12 \cdot 30 = (-4) \cdot 112 + 15 \cdot 30.$$

So in the case that $a = 112$ and $b = 30$, we have $d = 2$, $s = -4$, and $t = 15$.

We now summarise this section on Bézout's identity with a little scenario. Imagine a frog that lives on a doubly-infinite line of lilypads, indexed by the integers. It starts at the point 0 and can hop in steps of a or b (in either direction). Theorem 36 tells us that the lilypad $d = \gcd(a, b)$ is reachable by the frog. The next proposition characterises the set of *all* lilypads that the frog can get to.

Proposition 37. *Let a and b be nonzero integers. The set*

$$X = \{s'a + t'b : s', t' \in \mathbf{Z}\}$$

is exactly the set of multiples of $d = \gcd(a, b)$.

Proof. By Bézout's identity, there exist integers s and t such that $d = sa + tb$. First let $n \in \mathbf{Z}$ be a multiple of d . Then there is $k \in \mathbf{Z}$ such that $n = kd$, and we have

$$n = kd = d(sa + tb) = (ds)a + (dt)b,$$

which means that $n \in X$ (since ds and dt are both integers).

Conversely, suppose that $n \in X$, so $n = s'a + t'b$ for some $s', t' \in \mathbf{Z}$. Then since d divides a and d divides b , we can write $a = ld$ and $b = md$ for some integers $l, m \in \mathbf{Z}$. So we have

$$n = s'a + t'b = s'ld + t'md = (s'l + t'm)d,$$

which shows that $d \mid n$, since $s'l + t'm$ is an integer. \blacksquare

As a corollary, if $\gcd(a, b) = 1$, then it is possible for the robot to reach every integer! The situation in which $\gcd(a, b) = 1$ is very special, so much so that we have a name for it. We say that integers a and b are *relatively prime* or *coprime* if $\gcd(a, b) = 1$. By Bézout's identity, a and b are relatively prime if and only if there are integers s and t such that $1 = sa + tb$. This gives a very quick and easy way to check if certain numbers are relatively prime. For example, we have the following proposition.

Proposition 38. *For all integers $n > 1$, n and $n + 1$ are relatively prime.*

Proof. We have $1 = 1(n + 1) + (-1)n$. \blacksquare

9. Primes

An integer p is *prime* if $p \geq 2$ and for all $d \in \mathbf{N}$ with $d \mid p$, we either have $d = 1$ or $d = p$. An integer n is *composite* if $n \geq 2$ and n is not prime. (By negating the definition of prime, we see that $n \geq 2$ is composite if and only if there exist $a, b \in \{2, \dots, n-1\}$ such that $n = ab$.) Note that the integers 0 and 1 are neither prime nor composite.

The following theorem gives another characterisation of prime numbers.

Theorem 39. *An integer p with $p \geq 2$ is prime if and only if for all $a, b \in \mathbf{N}$, $p \mid ab$ implies that $p \mid a$ or $p \mid b$.*

Proof. First we prove the forward implication. Suppose that p is prime and let $a, b \in \mathbf{N}$ be such that $p \mid ab$. So there exists an integer k such that $ab = kp$. Consider $\gcd(a, p)$. Since the only divisors of p are 1 and p , this must be 1 or p . If it is p , then $p \mid a$ and we are done. So we restrict our attention to the case that $\gcd(a, p) = 1$, and our goal is to prove $p \mid b$. By Bézout's identity, there are integers s and t such that $1 = sa + tp$, so multiplying both sides by b , we have

$$b = bsa + btp = skp + btp = (sk + bt)p,$$

so $p \mid b$.

We prove the reverse implication by contraposition. Now suppose that p is composite. So there exist integers $2 \leq a, b \leq p-1$ such that $p = ab$. We want to show that p does not divide a and p does not divide b . We shall do this by showing that $u = a/p$ and $v = b/p$ are both not integers. Well since $p = ab = upb$, by dividing through by p we arrive at $1 = ub$, and since $b \geq 2$, $u = 1/b$ is between 0 and 1. This shows that p does not divide a . On the other

hand, since $p = ab = apv$, we have $1 = av$, and since $a \geq 2$, this means that $v = a/1$ is between 0 and 1. Hence p does not divide b . ■

To illustrate that p really does have to be prime for this theorem to hold, consider $p = 6$, $a = 2$ and $b = 15$. We have $6 \mid 30 = 2 \cdot 15$, but 6 divides neither 2 nor 15. By induction, the theorem extends to arbitrary finite products.

Corollary 40. *Let p be prime and n be a positive integer. If a_1, a_2, \dots, a_n are integers such that $p \mid a_1 a_2 \cdots a_n$, then $p \mid a_i$ for some $1 \leq i \leq n$.*

Proof. By induction on n . When $n = 1$, there is nothing to prove, for if $p \mid a_1$, then $p \mid a_i$ for $i = 1$.

Now suppose the statement holds for n . Assume that $p \mid a_1 a_2 \cdots a_n a_{n+1}$. By setting $a = a_1 a_2 \cdots a_n$ and $b = a_{n+1}$, we have $p \mid ab$, so by the previous theorem, either $p \mid a$ or $p \mid b$. If $p \mid b$, then $p \mid a_i$ for $i = n + 1$, and if $p \mid a$, then $p \mid a_1 \cdots a_n$, so by the induction hypothesis, $p \mid a_i$ for some $1 \leq i \leq n$. ■

Back in Section 6, we used the prime factorisation of integers in a proof, but didn't prove that statement itself. It's finally time to do so. Recall that we statement we used is that any integer n can be factored into a product

$$n = p_1^{v_1} p_2^{v_2} \cdots p_k^{v_k}$$

where $p_1 < p_2 < \cdots < p_k$ are primes and v_1, v_2, \dots, v_n are positive integers. By renumbering the primes, allowing them to possibly equal one another, we have the following equivalent statement.

Theorem 26' (*Fundamental Theorem of Arithmetic, again*). *Every integer $n \geq 2$ can be expressed as a product*

$$n = p_1 p_2 \cdots p_k$$

where $p_1 \leq p_2 \leq \cdots \leq p_k$ are prime numbers. Furthermore, this factorisation is unique.

Proof. The proof that such a decomposition of n exists is by (strong) induction. The base case is $n = 2$. This is already a prime factorisation, since 2 is prime.

For the inductive step, let $n \geq 2$ and assume such a decomposition exists for all $2 \leq i \leq n$. Now consider $n + 1$. If $n + 1$ is prime, then by setting $p_1 = n + 1$, we have a prime factorisation of $n + 1$ without getting out of bed. If $n + 1$ is not prime, then $n + 1 = ab$ for some integers $2 \leq a, b \leq n$. By the induction hypothesis, a and b can be factored into primes; that is, $a = p_1 \cdots p_l$ and $b = q_1 \cdots q_m$ for some primes $p_1, \dots, p_l, q_1, \dots, q_m$. So

$$n + 1 = ab = p_1 \cdots p_l q_1 \cdots q_m$$

is a factorisation of $n + 1$ into primes. It remains to arrange the p s and q s in nondecreasing order.

Now we prove that the prime factorisation of an integer is unique. In other words, $n \geq 2$ decomposes into

$$n = p_1 p_2 \cdots p_s$$

and

$$n = q_1 q_2 \cdots q_t$$

for primes $p_1, \dots, p_s, q_1, \dots, q_t$, then $s = t$ and $p_i = q_i$ for all $1 \leq i \leq s$. We also prove this statement by induction, but this time it is on the integer s . If $s = 1$, then

$$p_1 = n = q_1 q_2 \cdots q_t.$$

Note that t has to equal 1 here, since otherwise q_1 and q_2 would both divide p_1 and satisfy $2 \leq q_1, q_2 \leq p_1$, contradicting the fact that p_1 is prime. So $p_1 = q_1$ and we are done.

Next, suppose that the uniqueness statement holds for s (i.e., for any integer that decomposes into a product of s primes, the decomposition is unique) and we want to show that it holds for $s + 1$. We assume that

$$n = p_1 p_2 \cdots p_{s+1}$$

and

$$n = q_1 q_2 \cdots q_t,$$

where $p_1 \leq p_2 \leq \cdots \leq p_{s+1}$ and $q_1 \leq q_2 \leq \cdots \leq q_t$. We have $p_{s+1} \mid n = q_1 q_2 \cdots q_t$, so by the preceding corollary, there exists $1 \leq i \leq t$ such that $p_{s+1} \mid q_i$, and since q_i is prime, this means that $p_{s+1} = q_i \leq q_t$. Similarly, q_t divides $n = q_1 \cdots p_{s+1}$, so it divides p_j for some $1 \leq j \leq s + 1$. This means that $q_t = p_j \leq p_{s+1}$. Hence $q_t = p_{s+1}$. But

$$p_1 p_2 \cdots p_s p_{s+1} = q_1 q_2 \cdots q_t,$$

so we may divide out by $p_{s+1} = q_t$ on both sides to get

$$\frac{n}{p_{s+1}} = p_1 p_2 \cdots p_s = q_1 q_2 \cdots q_{t-1}.$$

By the induction hypothesis, the decomposition of n/p_{s+1} is unique, so $t - 1 = s$ and $p_i = q_i$ for all $1 \leq i \leq t - 1$. This proves that $s + 1 = t$ and $p_i = q_i$ for all $1 \leq i \leq s + 1$. ■

10.X We can use the Fundamental Theorem of Arithmetic to prove the following generalisation of Theorem 12.

Theorem 41. *Let k and n be positive integers. Then either $\sqrt[k]{n}$ is an integer or it is irrational.*

Proof. We prove this by contraposition, showing that if $\sqrt[k]{n}$ is rational, then it must be an integer. Suppose that $\sqrt[k]{n}$ is rational; so there exist integers a and b with $b \neq 0$ such that $\sqrt[k]{n} = a/b$. Without loss of generality, we may choose a and b such that $\gcd(a, b) = 1$. Writing

$$a = p_1^{v_1} p_2^{v_2} \cdots p_r^{v_r}$$

and

$$b = q_1^{w_1} q_2^{w_2} \cdots q_s^{w_s},$$

for primes $p_1 < p_2 < \cdots < p_r$, primes $q_1 < q_2 < \cdots < q_s$, and positive integers $v_1, v_2, \dots, v_r, w_1, w_2, \dots, w_s$, all the primes p_i and q_j must be different, since if there was a prime in common, call it $p = p_i = q_j$, we would have $p \mid \gcd(a, b)$, contradicting the fact that $\gcd(a, b) = 1$.

Taking the identity $\sqrt[k]{n} = a/b$ to the power of k yields

$$n = \frac{a^k}{b^k} = \frac{p_1^{kv_1} p_2^{kv_2} \cdots p_r^{kv_r}}{q_1^{kw_1} q_2^{kw_2} \cdots q_s^{kw_s}}.$$

Since there do not exist $1 \leq i \leq r$ and $1 \leq j \leq s$ such that $p_i = q_j$, this fraction is in reduced form. But since n is an integer, that means the denominator equals 1. In other words, $n = a^k$, so $\sqrt[k]{n} = a$ is an integer. ■

Let's see how to use this theorem. Suppose we want to know if the number $\sqrt[6]{18}$ is irrational or not. Well, since $1^6 = 1$ and $2^6 = 64$, we have $1^6 < 18 < 2^6$, meaning that $1 < \sqrt[6]{18} < 2$. (This is because the function $f : \mathbf{R} \rightarrow \mathbf{R}$ defined by $f(x) = x^6$ is increasing on the interval $[1, 2]$.) Hence $\sqrt[6]{18}$ is not an integer, and by the theorem it must be irrational.

Next, we show that there are infinitely many primes, a theorem first proved by Euclid.

Theorem 42. *There are infinitely many prime numbers.*

Proof. Suppose, towards a contradiction, that there are finitely many prime numbers, call them p_1, p_2, \dots, p_m . Let $n = p_1 p_2 \cdots p_m$, and consider the integer $n + 1$. Either it is prime or it is not. If $n + 1$ is prime, then we already have a contradiction, since $n + 1 > p_i$ for all $1 \leq i \leq m$, and we assumed that the p_1, \dots, p_m were all the primes. If $n + 1$ is not prime, then it is divisible by some prime in our list, call it p_i . Hence we can write $n = kp_i$ for some integer k , and we have

$$kp_i = n + 1 = p_1 p_2 \cdots p_m + 1.$$

Rearranging this a bit, we have

$$1 = kp_i - p_1 p_2 \cdots p_m = p_i(k - p_1 p_2 \cdots p_{i-1} p_{i+1} \cdots p_m);$$

in other words, p_i divides 1. This is a contradiction as 1 is not divisible by any integer greater than 1. ■

Though the set P of prime numbers is infinite, it does sort of get “sparser” as one heads off towards infinity. This is quantified by the following theorem, proved independently in 1896 by Jacques Hadamard and Charles Jean de la Vallée Poussin.

Theorem 43 (*Prime number theorem*). For $x \in \mathbf{R}$, let

$$\pi(x) = |\{p \leq x : p \text{ prime}\}|.$$

Then $\pi(x) \sim x/\ln x$, in the sense that

$$\lim_{x \rightarrow \infty} \frac{\pi(x)}{x/\ln x} = 1.$$

The proof is long and arduous, requiring a lot of background in complex analysis. It is often taught in a first-year graduate class in analytic number theory. Here is an easy corollary of the prime number theorem.

Corollary 44. Let n be a positive integer and let m be chosen uniformly at random from the set $\{1, 2, \dots, n\}$. Then

$$(\ln n) \mathbf{P}\{m \text{ prime}\} \rightarrow 1$$

as $n \rightarrow \infty$. In other words, the probability that m is prime is asymptotically $1/\ln n$.

Proof. Since m is chosen uniformly at random from $\{1, \dots, n\}$, the probability $\mathbf{P}\{m \text{ prime}\}$ equals $\pi(n)/n$. So

$$\lim_{n \rightarrow \infty} (\ln n) \mathbf{P}\{m \text{ prime}\} = \lim_{n \rightarrow \infty} \frac{\pi(n) \ln n}{n} = 1,$$

by the prime number theorem. ■

As an example, if we choose a 30-digit number at random (so $n = 10^{30} - 1$), the probability that this number is prime is roughly $1/\ln(10^{30}) = 1/(30 \ln 10) = 0.0145$ or 1.45%.

10. Modular arithmetic

Fix $n \geq 1$, and let $a, b \in \mathbf{Z}$. We say a is congruent to b modulo n if $n \mid a - b$, i.e., if $a - b = kn$ for some $k \in \mathbf{Z}$. Write $a \equiv b \pmod{n}$ or $a \equiv_n b$.

Take for example $n = 12$. We have, e.g., $4 \equiv 16 \pmod{12}$, since $4 - 16 = -12 = (-1)12$, but, e.g., $7 \not\equiv 17 \pmod{12}$, since $7 - 17 = -10$, and 12 does not divide -10 . This situation should be a familiar one, since we are used to working with numbers modulo 12 when telling the time.

For any fixed n , the set of all $(a, b) \in \mathbf{Z}^2$ with $a \equiv_n b$ is a relation on \mathbf{Z} . In fact, we have the following proposition.

Proposition 45. *For all fixed n , the relation \equiv_n is an equivalence relation on the set \mathbf{Z} .*

Proof. We must show that \equiv_n is reflexive, symmetric, and transitive.

Let $a \in \mathbf{Z}$. We have $a - a = 0 = 0 \cdot n$, so $a \equiv_n a$ (mod n), proving reflexivity.

Let $a, b \in \mathbf{Z}$ and suppose that $a \equiv_n b$, so there exists k such that $a - b = kn$. Then $b - a = (-k)n$, so $b \equiv_n a$. This proves symmetry.

Lastly, let $a, b, c \in \mathbf{Z}$ be such that $a \equiv_n b$ and $b \equiv_n c$. So there exist integers $k, l \in \mathbf{Z}$ such that $a - b = kn$ and $b - c = ln$. Then

$$a - c = (a - b) + (b - c) = kn + ln = (k + l)n,$$

which shows that $a \equiv_n c$. ■

Since \equiv_n is an equivalence relation, it partitions \mathbf{Z} into equivalence classes. We shall denote by $[a]_n$ the equivalence class of a modulo n . This is the set

$$[a]_n = \{b \in \mathbf{Z} : a - b = kn \text{ for some } k \in \mathbf{Z}\} = \{a + ln : l \in \mathbf{Z}\}.$$

For instance, when $n = 3$, the set

$$[0]_3 = \{\dots, -6, -3, 0, 3, 6, \dots\}$$

is just the set of all multiples of 3, and we also have

$$[1]_3 = \{\dots, -5, -2, 1, 4, 7, \dots\}$$

and

$$[2]_3 = \{\dots, -4, -1, 2, 5, 8, \dots\}.$$

These are all of the equivalence classes, since, for instance,

$$[4]_3 = \{\dots, -2, 1, 4, 7, 10, \dots\} = [1]_3.$$

Lastly, we shall touch upon the *modulo* operator, which is a feature of many programming languages. Let $a \in \mathbf{Z}$ and $b \geq 1$. Let $a \% b = r$ where q and r are the integers given by the division algorithm. (That is, $a = qb + r$ where $0 \leq r < b$.)

Proposition 46. *Fix an integer $n \geq 2$. Let $a, b \in \mathbf{Z}$. Then $a \equiv b \pmod{n}$ if and only if*

$$a \% n = b \% n.$$

Proof. Suppose that $a \equiv_n b$, so that there exists $k \in \mathbf{N}$ such that $a - b = kn$. Then, let $a = nq_1 + r_1$ and $b = nq_2 + r_2$ from the division algorithm. (So $r_1 = a \% n$ and $r_2 = b \% n$, and we have $0 \leq r_1, r_2 < n$.) Write

$$kn = a - b = (q_1 - q_2)n + (r_1 - r_2),$$

which we can rearrange to

$$r_1 - r_2 = kn - (q_1 - q_2)n = (k - q_1 + q_2)n.$$

This implies that $r_1 - r_2$ divides n , but since $r_1, r_2 \in [0, n)$, the quantity $r_1 - r_2$ is in the range $[-n + 1, n)$. Hence the only way it can divide n is for $r_1 - r_2 = 0$. We conclude that $a \% n = b \% n$.

For the other direction, we once again let $a = nq_1 + r_1$ and $b = nq_2 + r_2$ from the division algorithm. Now the assumption is that $r_1 = a \% n = b \% n = r_2$. So

$$a - b = (q_1 - q_2)n + (r_1 - r_2) = (q_1 - q_2)n.$$

This shows that $a \equiv b \pmod{n}$. ■

22.X This proposition is useful in practice. For example, suppose we want to know whether $22 \equiv 8 \pmod{3}$. We calculate $22 = 7 \cdot 3 + 1$, so $22 \% 3 = 1$; meanwhile $8 = 2 \cdot 3 + 2$, so $8 \% 3 = 2$. By the proposition, this means that $22 \not\equiv 8 \pmod{3}$.

The proposition also implies that for all $a \in \mathbf{Z}$ and $n \geq 2$ one has $[a]_n = [a \% n]$. From the division algorithm, we know that $a \% n$ is an element in the range $[0, n)$; it is equal to the integer r in that range such that we may write $a = qn + r$ for some integer q . We may choose to denote the whole equivalence class by this element r . The set

$$\mathbf{Z}/n\mathbf{Z} = \mathbf{Z}/\equiv_n = \{[0], [1], \dots, [n-1]\}$$

is called the *ring of integers modulo n* or the *cyclic group on n elements*.

Computations modulo n . We can do addition and multiplication modulo n by first doing addition and multiplication and then taking the remainder with respect to n . For example, when $n = 4$ we have the addition table

+	0	1	2	3
0	0	1	2	3
1	1	2	3	0
2	2	3	0	1
3	3	0	1	2

and the multiplication table

·	0	1	2	3
0	0	0	0	0
1	0	1	2	3
2	0	2	0	2
3	0	3	2	1

For $x, y \in \mathbf{Z}$, if $xy = 0$, then either $x = 0$ or $y = 0$. On the other hand, from the assumption $xy \equiv 0 \pmod{n}$ we cannot conclude in general that $x \equiv 0 \pmod{n}$ or $y \equiv 0 \pmod{n}$. As an example, when $n = 4$ and $x = y = 2$, neither x nor y is congruent to 0 modulo 4, but their product is 0 modulo 4. An element $a \in \mathbf{Z}$ with $a \not\equiv 0 \pmod{n}$ is said to be a *zero divisor* if there exists $b \in \mathbf{Z}$ with $b \not\equiv 0 \pmod{n}$ such that $ab \equiv 0 \pmod{n}$.

Remember that the numbers in the tables above are elements of $\mathbf{Z}/n\mathbf{Z}$ (in the case $n = 4$), hence they are not really integers so much as equivalence classes of integers, which are sets! So it is not immediate that addition and multiplication are well-defined modulo some fixed integer n . For instance, to say that $[3]_4 \cdot [1]_4 = [3]_4$, it is not enough to just observe that $3 \cdot 1 = 3$. We need to show that for any integers $a \in [3]_4$ and $b \in [1]_4$, we have $a \cdot b \in [3]_4$. This is the essence of the following proposition.

Proposition 47. *Let $a \equiv c \pmod{n}$ and $b \equiv d \pmod{n}$. Then*

- i) $a + b \equiv c + d \pmod{n}$;
- ii) $ab \equiv cd \pmod{n}$; and
- iii) $a^m \equiv c^m \pmod{n}$ for all $m \in \mathbf{N}$.

Proof. By assumption, there are integers k and l such that $a - c = kn$ and $b - d = ln$.

To prove (i), we compute

$$(a + b) - (c + d) = a - c + b - d = kn + ln = (k + l)n,$$

so $a + b \equiv c + d \pmod{n}$.

For part (ii), we compute

$$\begin{aligned} ab - cd &= ab - cb + cb - cd \\ &= b(a - c) + c(b - d) \\ &= bkn + cln \\ &= (bk + cl)n, \end{aligned}$$

so $ab \equiv cd \pmod{n}$.

We prove part (iii) by induction. The base case is $m = 0$; we have $a^0 = 1 = c^0$, so $a^0 \equiv c^0 \pmod{n}$. Now let $m \geq 0$ and assume that $a^m \equiv c^m \pmod{n}$. Combining this with the hypothesis $a \equiv c \pmod{n}$, we use part (ii) to conclude that $a \cdot a^m \equiv c \cdot c^m \pmod{n}$, which simplifies to $a^{m+1} \equiv c^{m+1} \pmod{n}$. \blacksquare

We can use this to perform the modulo operation on large numbers without necessarily knowing the large numbers themselves. For example, let $m = 21^{10} - 3 \cdot 13 + 14$ and suppose we want to know $m \% 6$.

First we observe that $14 \equiv 2 \pmod{6}$ and $13 \equiv 1 \pmod{6}$. So

$$m = 21^{10} - 3 \cdot 13 + 14 \equiv_6 21^{10} - 3 \cdot 1 + 2 \equiv_6 21^{10} - 1.$$

Now we deal with the exponent by first taking the base of the exponent modulo 6. We have $21 \equiv 3 \pmod{6}$; hence

$$m \equiv_6 3^{10} - 1.$$

To reduce the exponent, we factorise $10 = 2 \cdot 5$, yielding

$$m \equiv_6 (3^2)^5 - 1.$$

But again we can reduce the base here, since $9 \equiv_6 3$, which gives us

$$m \equiv_6 3^5 - 1.$$

Now since the exponent 5 is prime, we aren't able to reduce it by factorising. However, we do have $5 = 2+2+1$, and, using once again the fact that $3^2 = 9 \equiv_6 3$, we obtain

$$m \equiv_6 3^2 \cdot 3^2 \cdot 3 - 1 \equiv_6 3 \cdot 3 \cdot 3 - 1 \equiv_6 3 \cdot 3 - 1 \equiv_6 3 - 1 \equiv_6 2.$$

Inverses modulo n . Considering 2 as a rational number, we say that the inverse of 2 is $1/2$, since $2 \cdot (1/2) = 1$. We write $2^{-1} = 1/2$. We shall use the same terminology in the ring of integers modulo n . An element $a \in \mathbf{Z}$ is said to be *invertible modulo n* if there exists $b \in \mathbf{Z}$ such that $ab \equiv 1 \pmod{n}$. In this case we shall say that b is an *inverse* of a . In fact, inverses are unique (in $\mathbf{Z}/n\mathbf{Z}$).

Proposition 48. *Let $a, n \in \mathbf{Z}$ with $n \geq 2$. Then if $ab \equiv 1 \pmod{n}$ and $ac \equiv 1 \pmod{n}$, then $b \equiv c \pmod{n}$.*

Proof. We have

$$b \equiv b \cdot 1 \equiv b(ac) \equiv (ab)c \equiv 1 \cdot c \equiv c \pmod{n}. \quad \blacksquare$$

Hence we may speak of *the* inverse of an element a , which we shall denote by a^{-1} . We have $a^{-1}a \equiv 1 \equiv aa^{-1} \pmod{n}$.

For instance, when $n = 5$, the fact that

$$2 \cdot 3 \equiv_5 6 \equiv_5 1$$

tells us that 2 (and 3, for that matter) are invertible modulo 5. When $n = 6$, is 2 still invertible? Here n is small enough that we can just try all the cases. We have $2 \cdot 0 = 0$, $2 \cdot 1 = 2$, $2 \cdot 2 = 4$, $2 \cdot 3 = 6 \equiv_6 0$, $2 \cdot 4 = 8 \equiv_6 2$, and $2 \cdot 5 = 10 \equiv_6 4$. Since none of these were equivalent to 1 modulo 6, we see that 2 is not invertible modulo 6.

It was quite a tedious process to exhaustively show that 2 is not invertible modulo 6. The following theorem gives a much easier-to-use criterion for invertibility.

Theorem 49. Let $a, n \in \mathbf{Z}$ with $n \geq 2$. Then

- i) a is invertible modulo n if and only if $\gcd(a, n) = 1$; and
- ii) if a is invertible, then there is a unique integer $b \in [0, n-1]$ such that $ab \equiv 1 \pmod{n}$. Namely, if

$$1 = sa + tn,$$

then we can set $b = s \% n$.

Proof. The proof of (i) is a fairly simple corollary of Bézout's theorem. If a is invertible, then $ab \equiv 1 \pmod{n}$, so there exists $k \in \mathbf{Z}$ such that $1 - ab = kn$. Rearranging this, we have $1 = ab + kn$, and by Bézout's theorem (specifically, the final statement of Theorem 36, this implies that $\gcd(a, n) = 1$. On the other hand, if $\gcd(a, n) = 1$, then there exists $k \in \mathbf{Z}$ such that $1 = ab + kn$. This means that $1 - ab = kn$, so $ab \equiv 1 \pmod{n}$.

For part (ii), suppose that $s, t \in \mathbf{Z}$ are such that $1 = sa + tn$. Then taking this equation modulo n , we have

$$1 = sa + tn \equiv_n sa + 0t \equiv_n sa.$$

Letting $b = s \% n$, we have $1 = ba$, since s and b belong to the same equivalence class modulo n , and b is an integer in $[0, n-1]$ by the division algorithm. This proves existence of such a b . To show uniqueness, we suppose that b and c are both elements of $[0, n-1]$ with $ab \equiv 1 \pmod{n}$ and $ac \equiv 1 \pmod{n}$. By Proposition 48, we have $b \equiv c \pmod{n}$; in other words, $n \mid b - c$. But since $b, c \in [0, n-1]$, we have $b - c \in [-n+1, n-1]$, and the only divisor of n in this range is 0. Hence $b - c = 0$ and we conclude that $b = c$. ■

We can use this theorem to systematically find inverses of integers modulo other integers. For instance, suppose we want to find the inverse of 17 modulo 20. First we find $\gcd(20, 17)$ by Euclid's algorithm:

$$\begin{aligned} 20 &= 1 \cdot 17 + 3 \\ 17 &= 5 \cdot 3 + 2 \\ 3 &= 1 \cdot 2 + 1 \\ 1 &= 1 \cdot 1 + 0 \end{aligned}$$

The fact that $\gcd(20, 17) = 1$ tells us that 17 indeed has an inverse. Now we go backwards to express 1 as an integer linear combination of 20 and 17:

$$\begin{aligned} 1 &= 3 - 1 \cdot 2 \\ &= 3 - (17 - 5 \cdot 3) \\ &= -17 + 6 \cdot 3 \\ &= -17 + 6(20 - 17) \\ &= 6 \cdot 20 - 7 \cdot 17 \end{aligned}$$

So the inverse of 17 modulo 20 is $-7 \equiv_{20} 13$. (We can verify that

$$17 \cdot 13 \equiv (-3)(-7) \equiv 21 \equiv 1 \pmod{20}.$$

As a trick for computations, use the range $-n/2$ to $n/2$ instead of 0 to $n-1$ to keep the numbers a bit smaller. It might help you avoid errors.)

Let's find the full list of invertible integers modulo 20, as well as their inverses. First, we list the integers in the range $[0, 19]$ that are relatively prime with 20:

$$1, 3, 7, 9, 11, 13, 17, 19$$

If some number is in this list, then its inverse must also be in this list (since inverses to a given element must themselves be invertible). The elements 1 and 19 are their own inverses, since $1 \cdot 1 = 1$ and $19 \cdot 19 \equiv_{20} (-1)(-1) = 1$. We already saw above that 13 and 17 are inverses to each other. That leaves 3, 7, 9, and 11. First we see that $3 \cdot 7 = 21 \equiv_{20} 1$, so 3 and 7 are inverses to each other. Then we note that $9 \cdot 9 = 81 \equiv_{20} 1$, so 9 is its own inverse, leaving us to conclude that 11 must also be its own inverse as well—either by the fact that we've ruled out all other possibilities, or by the computation

$$11 \cdot 11 \equiv_{20} (-9)(-9) = 81 \equiv_{20} 1.$$

We now investigate what happens in the particular case where the modulus n is prime. In this situation we have the following proposition.

Proposition 50. *Let p be prime. Then*

- i) *every $x \in \mathbf{Z}$ with $x \not\equiv 0 \pmod{p}$ is invertible modulo p ; and*
- ii) *for all $a, b \in \mathbf{Z}$ with $ab \equiv 0 \pmod{p}$ one has $a \equiv 0 \pmod{p}$ or $b \equiv 0 \pmod{p}$.*

Proof. Since p only has factors 1 and p , $\gcd(x, p)$ must either be 1 or p . But the condition $x \not\equiv 0 \pmod{p}$ implies that p does not divide x . Hence $\gcd(x, p) = 1$, proving part (i).

For part (ii), assume that $ab \equiv 0 \pmod{p}$, so $p \mid ab$. But p is prime, so by Theorem 39, either $p \mid a$ or $p \mid b$. In the first case we have $a \equiv 0 \pmod{p}$, and in the second case, $b \equiv 0 \pmod{p}$. ■

24.X **Solving congruences modulo n .** Suppose we want to solve for all integers x satisfying $x^2 \equiv x \pmod{n}$, first for $n = 6$, then for $n = 7$.

For the case where $n = 6$, it suffices to consider $x \in \{0, 1, \dots, n-1\}$, since if $x^2 \equiv x \pmod{n}$ and $x \equiv a \pmod{n}$, then $a^2 \equiv a \pmod{n}$. So we simply try all $x \in \{0, 1, \dots, 5\}$. We have $0^2 = 0$, $1^2 = 1$, $2^2 = 4 \not\equiv_6 2$, $3^2 = 9 \equiv_6 3$, $4^2 = 16 \equiv_6 4$, and $5^2 \equiv_2 (-1)^2 = 1 \not\equiv_6 5$. We conclude that 0, 1, 3, and 4 are solutions to this congruence modulo 6, and more broadly, any integer y of the form

$$y = x + 6k$$

where $k \in \mathbf{Z}$ and $x \in \{0, 1, 3, 4\}$, is a solution $x^2 \equiv x \pmod{6}$.

Now we tackle the case where $n = 7$. In fact, the solution would be no different should n be any other prime. First we subtract x from both sides, obtaining the congruence $x^2 - x \equiv 0 \pmod{7}$, then we factorise the left side to get $x(x - 1) \equiv 0 \pmod{7}$. Now by the previous proposition, since 7 is prime if $x(x - 1) \equiv 0 \pmod{7}$, we must have either $x \equiv 0 \pmod{7}$ or $x - 1 \equiv 0 \pmod{7}$, so the only solutions are $x \equiv_7 0$ and $x \equiv_7 1$ (and anything in their equivalence classes).

Replacing 7 with any other prime p in the above paragraph yields a proof of the following proposition.

Proposition 51. *Let p be a prime. Then $a^2 \equiv a \pmod{p}$ if and only if a is either congruent to 0 or 1 modulo p . ■*

Here is a similar proposition.

Proposition 52. *Let p be a prime and let $a \not\equiv 0 \pmod{p}$ (so that a is invertible). Then $a \equiv a^{-1} \pmod{p}$ if and only if a is either congruent to 1 or -1 modulo p .*

Proof. If $a = 1$ or $a = -1$, then $a^2 = 1$, so $a^2 \equiv 1$, and multiplying by a^{-1} on both sides, we have $a \equiv a^{-1} \pmod{p}$.

Now suppose that $a \equiv a^{-1} \pmod{p}$. Multiplying by a on both sides, we get $a^2 \equiv 1 \pmod{p}$; then, subtracting 1 from both sides and factoring the resulting polynomial, we get

$$(a + 1)(a - 1) \equiv 0 \pmod{p}.$$

This implies that either $a + 1 \equiv 0 \pmod{p}$ or $a - 1 \equiv 0 \pmod{p}$. In the first case, $a \equiv_p -1$, and in the second case, $a \equiv_p 1$. ■

We finish off this section with an important theorem, first stated by Pierre de Fermat in 1640. He did not supply a proof; the first published proof of this theorem was given by Leonhard Euler in 1736.

Theorem 53 (*Fermat's little theorem*). *Let a and p be integers with p prime. If $a \not\equiv 0 \pmod{p}$, then $a^{p-1} \equiv 1 \pmod{p}$.*

Before proving this theorem, we first state and prove a lemma. Recall that we define the *factorial* of $n \in \mathbf{N}$ to be the product $n! = 1 \cdot 2 \cdots (n - 1)n$.

Lemma 54. *For all prime numbers p , the integer $(p - 1)!$ is congruent to -1 modulo p .*

Proof. If $p = 2$, we have $(p - 1)! = 1 \equiv -1 \pmod{2}$, and if $p = 3$, then $(p - 1)! = 2 \equiv -1 \pmod{3}$.

For the rest of the proof, assume that $p \geq 5$. By Proposition 52, each element in the set $S = \{2, 3, \dots, p - 3, p - 2\}$ is not its own inverse modulo p , so for each element $s \in S$, there is some other element $s' \in S$ with $s' \neq s$ such that $ss' \equiv 1 \pmod{p}$. This means that

$$2 \cdot 3 \cdots (p - 3)(p - 2) \equiv 1 \pmod{p}.$$

From this, we see that

$$(p-1)! = 1 \cdot 2 \cdots (p-2)(p-1) = (2 \cdot 3 \cdots (p-3)(p-2))(p-1) \equiv_p (p-1) \equiv_p -1. \quad \blacksquare$$

We are now able to prove Fermat's little theorem.

Proof of Theorem 53. Since $a \not\equiv 0 \pmod{p}$, it is invertible modulo p by Proposition 50. Denote its inverse by a^{-1} . Let $G = (\mathbf{Z}/p\mathbf{Z}) \setminus \{0\}$ for short. Define a function $f : G \rightarrow G$ by letting $f(x) = a \cdot x$, where we consider the result modulo p (and hence the result is an element of $\mathbf{Z}/p\mathbf{Z}$). This function is well defined: since a and x are both not zero modulo p , their product will be an element of G . If $f(x_1) = f(x_2)$, then

$$a \cdot x_1 \equiv a \cdot x_2 \pmod{p},$$

so multiplying both sides by a^{-1} , we have $x_1 \equiv x_2 \pmod{p}$. This proves that f is injective. Now let $y \in G$. We want to find x with $f(x) \equiv y \pmod{p}$. To do so, simply set $x \equiv a^{-1}y \pmod{p}$. Then

$$f(x) = f(a^{-1}y) = a(a^{-1}y) \equiv y \pmod{p}.$$

This proves that f is surjective, and hence bijective.

The bijection f shows that the set

$$\{a, 2a, \dots, (p-2)a, (p-1)a\},$$

taken as a subset of $\mathbf{Z}/p\mathbf{Z}$ (i.e., we take each element modulo p , in the range $[0, \dots, p-1]$), is exactly the same as the set

$$\{1, 2, \dots, p-2, p-1\},$$

just that the order of elements might be permuted. Hence we have

$$a(2a) \cdots ((p-2)a)((p-1)a) \equiv (p-1)! \pmod{p}.$$

Combining all the a factors on the left-hand side, we get

$$a^{p-1}(p-1)! \equiv (p-1)! \pmod{p}.$$

But by the previous lemma, -1 is the inverse of $(p-1)!$ modulo p , so multiplying both sides of this congruence by -1 , we get

$$a^{p-1} \equiv 1 \pmod{p},$$

which is what we wanted to show. \blacksquare

As a matter of interest, the intermediary lemma we proved is one direction of Wilson's theorem, proved by John Wilson in 1770.

Theorem 55 (*Wilson's theorem*). For all integers $n \geq 2$, the congruence

$$(n-1)! \equiv -1 \pmod{n}$$

holds if and only if n is prime.

Proof. We already proved the “if” direction earlier, as Lemma 54. We leave the “only if” direction of the proof as an exercise for the reader. ■

*Tout nombre premier mesure infalliblement
une des puissances - 1 de quelque progression que ce soit,
& l'exposant de ladite puissance est sous-multiple du nombre premier donné - 1.
Et après qu'on a trouvé la première puissance qui satisfait à la question,
toutes celles dont les exposants sont multiples de l'exposant de la première
satisfont de même à la question.*

— PIERRE DE FERMAT, in a letter to Bernard Frénicle de Bessy (1640)

11. Applications of number theory

In this section we present a potpourri of interesting ways number theory is applied to make your life better.

ISBN book identifiers. Every published book has an ISBN code that serves as its unique identifier. This code is 10 digits long for books published before 2007, but three new digits have been added for books published after 2007. In this section we'll deal with the simpler case of 10 digits.

What we want to do is to bake some redundancy into the codes, so that if a single digit is typed wrong, a computer will be able to tell the user that the code is invalid. This allows the user to correct their mistake. The way we do this, in a ten-digit code $d_{10}d_9d_8 \cdots d_2d_1$, is to have only the first nine digits encode the book. The last digit is a *check digit*, chosen so that

$$\sum_{i=1}^{10} id_i \equiv 0 \pmod{11}. \quad (4)$$

In other words,

$$d_1 \equiv -\sum_{i=2}^{10} id_i \pmod{11}.$$

If d_1 needs to equal 10, we use the symbol ‘X’.

How does this solve our problem? Well, if a hapless person takes a valid ISBN code and mangles it by either getting one digit wrong, or by swapping two adjacent digits, the result is an invalid ISBN code, and the computer will be able to flag it. We formalise this statement as the following theorem.

Theorem 56. Let $d_{10}d_9d_8 \cdots d_2d_1$ be a valid ISBN code; that is, a code that satisfies (4). Then any code obtained by either

- i) changing exactly one digit d_i , for some $1 \leq i \leq 10$, or
- ii) swapping distinct adjacent digits d_i and d_{i-1} , for some $2 \leq i \leq 10$,

is not a valid ISBN codes (it does not satisfy (4)).

Proof. First we prove (i). Suppose that $c_{10}c_9c_8 \cdots c_2c_1$ is a code that is equal to $d_{10}d_9d_8 \cdots d_2d_1$ except at one place $1 \leq j \leq 10$. So $c_j \neq d_j$, but for all $1 \leq i \leq 10$ with $i \neq j$, we have $c_i = d_i$. Then

$$\begin{aligned} \sum_{i=1}^{10} ic_i &\equiv_{11} \sum_{i=1}^{10} ic_i - 0 \\ &\equiv_{11} \sum_{i=1}^{10} ic_i - \sum_{i=1}^{10} id_i \\ &= \sum_{i=1}^{10} i(c_i - d_i) \\ &= j(c_j - d_j). \end{aligned}$$

But $j \in \{1, \dots, 10\}$ and $(c_j - d_j) \in [-10, 10]$. So neither j nor $c_j - d_j$ are congruent to 0 modulo 11. Since 11 is prime, their product cannot be congruent to 0 modulo 11. We conclude that $\sum_{i=1}^{10} ic_i \not\equiv 0 \pmod{11}$; i.e., this is not a valid ISBN code.

To prove part (ii), suppose that $c_{10}c_9c_8 \cdots c_2c_1$ is a code that is equal to $d_{10}d_9d_8 \cdots d_2d_1$, except that for some $2 \leq j \leq 10$, we have $c_j = d_{j-1}$ and $c_{j-1} = d_j$ (the adjacent digits are swapped). Furthermore, assume that $d_j \neq d_{j-1}$, since swapping adjacent digits that are the same doesn't do anything to the code. Then

$$\begin{aligned} \sum_{i=1}^{10} ic_i &\equiv_{11} \sum_{i=1}^{10} ic_i - 0 \\ &\equiv_{11} \sum_{i=1}^{10} ic_i - \sum_{i=1}^{10} id_i \\ &= \sum_{i=1}^{10} i(c_i - d_i) \\ &= (j-1)c_{j-1} + jc_j - (j-1)d_{j-1} - jd_j \\ &= (j-1)d_j + jd_{j-1} - (j-1)d_{j-1} - jd_j \\ &= d_{j-1} - d_j. \end{aligned}$$

But since d_{j-1} and d_j are both in the range $\{0, \dots, 10\}$, their difference is in the range $[-10, 10]$, and we assumed that $d_{j-1} \neq d_j$, so this difference is not zero. We conclude that $d_{j-1} - d_j \not\equiv 0 \pmod{11}$, so the code is not valid. ■

Something similar is done with credit card numbers, though we don't take digits modulo 11 (or else we'd have to use the digit 'X', which is cumbersome). Instead, a slightly different method called Luhn's algorithm is used. The calculation used to find the check digit is a bit different (let's not get into it), but the result is that we can still detect single-digit errors, and we can detect all swapping errors except for the transposition $09 \leftrightarrow 90$.

Divisibility tests. In grade school you might have learned the following divisibility rule: *To tell if a number is divisible by 3, you sum up all of its digits. If this is a multiple of 3, then the original number was a multiple of 3, and if not, then the original number wasn't.* The same holds with 3 replaced by 9.

Using the number theory we've learned, we are now able to prove this fact.

Proposition 57. *Suppose that $n \in \mathbf{N}$ is written in base-10 digits as*

$$n = d_k d_{k-1} d_{k-2} \cdots d_2 d_1 d_0,$$

for some integer $k \geq 0$ and $d_i \in \{0, \dots, 9\}$ for $0 \leq i \leq k$. Then n is divisible by 3 if and only if

$$\sum_{i=0}^k d_i$$

is also divisible by 3.

Proof. Since $10 \equiv 1 \pmod{3}$, any power of 10 is also congruent to 1 modulo 3. So we have

$$n = \sum_{i=0}^k d_i 10^i \equiv_3 \sum_{i=0}^k d_i,$$

and we see that n is divisible by 3 if and only if the sum on the right-hand side is. ■

This works with 3 replaced by 9 since we also have $10 \equiv 1 \pmod{9}$.

Computing large powers modulo n . The next application is not a precise mathematical statement; rather, it's an observation that Fermat's little theorem can allow us to compute what large powers are modulo a prime without having to compute the number itself. We'll illustrate this by two examples.

First we compute $25^{134} \pmod{11}$. We begin by noting that

$$25 \equiv 3 \pmod{11},$$

so what we really want is $3^{134} \pmod{11}$. By Fermat's little theorem, $3^{10} \equiv 1 \pmod{11}$, so we may cast out multiples of 10 from 134 to arrive at the easy computation

$$3^{134} \equiv 3^{10 \cdot 13 + 4} \equiv (3^{10})^{13} \cdot 3^4 \equiv 3^2 \cdot 3^2 \equiv 9 \cdot 9 \equiv (-2)(-2) \equiv 4 \pmod{11}.$$

Now let's try $25^{134} \pmod{14}$. Since 14 is not prime, we can no longer use Fermat's little theorem. The trick we shall use is to square 25 repeatedly. First note that $25 \equiv -3 \pmod{14}$, so $25^2 \equiv 9 \pmod{14}$. Then, squaring again, we get

$$25^4 \equiv 9^2 \equiv (-5)^2 \equiv 25 \equiv -3 \pmod{14}.$$

This means that repeated squaring of 25 cycles between -3 and 9 modulo 14, so we immediately have

$$25^8 \equiv 9, \quad 25^{16} \equiv -3, \quad 25^{32} \equiv 9, \quad 25^{64} \equiv -3,$$

and

$$25^{128} \equiv 9$$

modulo 14. Now writing the exponent 134 as a sum of powers of two allows us to conclude that

$$25^{134} = 25^{128+4+2} = 25^{128} \cdot 25^4 \cdot 25^2 \equiv 9 \cdot (-3) \cdot 9 \equiv (-3)(-3) \equiv 9 \pmod{14}.$$

Fermat's primality test. Fermat's little theorem says that if n is prime, then for all integers a with $a \not\equiv 0 \pmod{n}$ one has $a^{n-1} \equiv 1 \pmod{n}$. The contrapositive of this statement is that if there is some integer a with $a \not\equiv 0 \pmod{n}$ and $a^{n-1} \not\equiv 1 \pmod{n}$, then n is not prime. This leads us to the following algorithm that can (usually) tell us when a number is not prime.

Algorithm F (*Fermat's primality test*). Given an integer $n \geq 4$, this algorithm will attempt to declare that n is not prime. If the algorithm fails to do so, then the test is inconclusive (n may either be prime or not).

F1. [Initialise.] Set $a \leftarrow 2$. (We can skip the step $a = 1$ since a^{n-1} will always equal 1 in this case.)

F2. [Compute power.] Set $b \leftarrow a^{n-1} \% n$. (So $0 \leq b < n$.)

F3. [Not prime?] If $b \neq 1$, output "Not prime," and terminate the algorithm.

F4. [Loop?] Set $a \leftarrow a + 1$. If $a = n$, output "Inconclusive," and terminate the algorithm. Otherwise, go to step F2. ■

29.X Let's see a couple of examples. Suppose we want to find out if 9 is prime or not (we both know it isn't but play along for a second). In the very first step loop, when $a = 2$, we have

$$2^{9-1} = 2^8 = (2^4)^2 = 16^2 \equiv (-2)^2 = 4 \pmod{9},$$

so since $4 \not\equiv 1 \pmod{9}$, we conclude right away that 9 is not prime.

A harder one. Is 341 prime? We start with $a = 2$. Noting that $3 \cdot 341 = 1023$, we have

$$2^{340} = (2^{10})^{34} = 1024^{34} \equiv 1^{34} \equiv 1 \pmod{341},$$

so the test is inconclusive after the first loop. Fine. If we try $a = 3$ next, we will find that $3^{340} \equiv 56 \not\equiv 1 \pmod{341}$, so we conclude that 341 is not prime. (Those wishing to practise computing large powers modulo n might like to verify the fact that $3^{340} \equiv 56 \pmod{341}$. Use the squaring trick.)

If n causes Algorithm F to output “Inconclusive,” it is very likely that n is prime. However, there are rare cases of composite numbers n that cause Algorithm F to output “Inconclusive.” These are called *Carmichael numbers*. The smallest one is 561.

Vernam’s one-time pad. Alice and Bob want to exchange secret messages over a public channel without anyone being able to decipher them. A message can be encoded as a binary string M . Here is one way that Alice and Bob can securely send encrypted messages to each other. It is called Vernam’s one-time pad, named for G. Vernam, an AT&T engineer who patented the method in 1917.

First, they need to meet in person and generate a long sequence of random bits. (How can one generate random numbers? This is an interesting quasi-philosophical question.) Each of them keeps a copy of this sequence, call it S .

Now suppose Alice and Bob separate and Alice wants to send the message $M = 0110$ to Bob. Of course, she shouldn’t send this string directly, as it could be intercepted by foes. She needs to use the sequence S . Suppose the first four bits of S are 1101. To encrypt her message, Alice adds each bit of M to the corresponding bit of S , modulo 2. We will denote this by the \oplus operation. The encrypted message is $M' = 0110 \oplus 1101 = 1011$. She sends it to Bob.

Bob receives the string $M' = 1011$. To decode it and recover the message M , Bob simply has to perform the \oplus operation using the (same) first four bits of S that Alice used to encode the message in the first place. Indeed, we see that $1101 \oplus 1101 = 0110 = M$. The reason that this works is that for any binary string S , we have $S \oplus S = 0$, so performing the \oplus operation twice always returns the original string.

If they want to send more messages, they simply discard the bits they used before and continue using the next few bits of S .

The advantage to this approach is that an adversarial third-party, reading M' , gains no information about the original string M , provided they do not know the sequence S . In other words, it is entirely secure. The disadvantages of Vernam’s one-time pad are that Alice and Bob must meet in person beforehand, they need to find a way to generate random bits securely, and they need to periodically repeat this process to replenish the bits once they run out.

The RSA cryptosystem. Now we move on to a more sophisticated method of encrypting and decrypting messages. It is called the RSA *cryptosystem*, named for R. Rivest, A. Shamir, and L. Adleman, who described the algorithm in 1977.

Here’s the scenario. Alice wants people to be able to send her messages such that only she can decode them. To do so, she needs to create a *public key*, which allows people to encrypt their messages before sending them to her, and then a

private key that will allow her to decrypt messages. (Only she knows the private key.)

To set things up, Alice needs two large primes, call them p and q . (Something like 200 bits is sufficient.) Then she computes $n = pq$. Now, she picks some integer k with

$$\gcd(k, (p-1)(q-1)) = 1.$$

By Theorem 49, there is some integer s such that $ks \equiv 1 \pmod{(p-1)(q-1)}$. Alice publishes the integer n and the public key k . The private key is s ; this she keeps to herself.

Now Bob wants to send a secret message M to Alice. Considering M as an integer (in binary), a requirement for the algorithm to work is $M < n$. Bob computes $\overline{M} = M^k \% n$ and sends \overline{M} to A. (We have $0 \leq \overline{M} < n$ by the division algorithm.)

Alice receives \overline{M} . She computes \overline{M}^s . It turns out that this is M . To prove this, we need a lemma.

Lemma 58. *Let a and b be integers.*

- i) *For all positive integers m and n , if $a \equiv b \pmod{mn}$, then $a \equiv b \pmod{m}$ and $a \equiv b \pmod{n}$.*

In the case that m and n are different primes, the converse holds; that is, the following holds.

- ii) *For all primes p and q with $p \neq q$, if $a \equiv b \pmod{p}$ and $a \equiv b \pmod{q}$, then $a \equiv b \pmod{pq}$.*

Proof. Assume that $a \equiv b \pmod{mn}$, so there exists $k \in \mathbf{Z}$ such that $a - b = kmn$. This means that $a - b \equiv 0 \pmod{m}$ and $a - b \equiv 0 \pmod{n}$, so $a \equiv b \pmod{m}$ and $a \equiv b \pmod{n}$.

To prove (ii), suppose that there exist $k, l \in \mathbf{Z}$ such that $a - b = kp$ and $a - b = lq$. We have $kp = lq$, so $p \mid lq$. Because p is prime, p must divide either l or q . But q is prime as well, so the only divisors of q are 1 and q . Hence $p \mid l$, and we may write $l = rp$ for some $r \in \mathbf{Z}$. So we have

$$a - b = lq = rpq,$$

implying that $a - b \equiv 0 \pmod{pq}$; that is, $a \equiv b \pmod{pq}$. ■

Now we shall prove that the encrypted message $\overline{M} = M^k$ can be decrypted into the original message.

Theorem 59 (RSA encryption). *Let p and q be distinct primes and let $n = pq$. Let k and s be such that $ks \equiv 1 \pmod{n}$. Then for all integers $0 \leq M < n$, $(M^k)^s \equiv M \pmod{n}$.*

Proof. Without loss of generality, we can take k and s both positive. From $ks \equiv 1 \pmod{(p-1)(q-1)}$, we know there is an integer l such that $ks - 1 = l(p-1)(q-1)$.

Since $k, s > 0$, this identity means that $l > 0$ as well. Now, since $n = pq$, we have

$$(M^k)^s \equiv M^{ks} \equiv M^{1+l(p-1)(q-1)} \equiv M(M^{p-1})^{l(q-1)} \pmod{pq}.$$

By part (i) of the previous lemma, this means we have the same congruence modulo p ; i.e.,

$$(M^k)^s \equiv M(M^{p-1})^{l(q-1)} \pmod{p},$$

and by Fermat's little theorem $M^{p-1} \equiv 1$ for all $M \not\equiv 0 \pmod{p}$, so this simplifies to

$$\begin{aligned} (M^k)^s &\equiv_p \begin{cases} M, & \text{if } M \not\equiv 0 \pmod{p}; \\ 0, & \text{otherwise} \end{cases} \\ &\equiv_p M. \end{aligned}$$

Mutatis mutandis (by the roles of p and q in the above), we have $(M^k)^s \equiv_q M$ as well, so part (ii) of the previous lemma can now be applied to give

$$(M^k)^s \equiv M \pmod{pq}. \quad \blacksquare$$

Here is an example of the RSA algorithm in action. Alice sets $p = 3$ and $q = 11$, so that $n = 33$. This means that $(p-1)(q-1) = 2 \cdot 10 = 20$. Using $k = 7$ (which we can do since $\gcd(7, 20) = 1$, we have $s = 3$, since $3 \cdot 7 = 21 \equiv 1 \pmod{20}$). Alice publishes the integer $n = 33$ as well as the public key $k = 7$. She keeps the number $s = 3$ private.

The possible messages M one can send to Alice are $0, 1, \dots, 32$, since we need $0 \leq M < n$. Suppose Bob wants to send the message 00010 to Alice; that is, we have $M = 2$. Bob computes

$$M^k = 2^7 = 2^5 \cdot 2^2 = 32 \cdot 4 \equiv (-1) \cdot 4 \equiv 29 \pmod{33},$$

and sends Alice $\overline{M} = 29$.

Now Alice computes

$$\overline{M}^s = 29^3 \equiv (-4)^3 = -64 \equiv -64 + 66 = 2 \pmod{33}$$

and recovers the original message $M = 2$.

Now imagine that there is a third person, Charlie, who wants to thwart Alice and Bob. The integer n and the public key n are known to all, including Charlie. If he wants to decrypt Bob's message, all Charlie needs to know is the the private key s . In this case, since n is very small, the Adversary can simply factor $33 = 3 \cdot 11$ to find out that $(p-1)(q-1) = 20$, and then use Euclid's Algorithm to find the inverse of 7 modulo 20, and he will be able to decrypt all messages between Alice and Bob.

In real life, RSA is implemented using very large primes p and q . This makes it difficult, but not impossible, to break the encryption. To put things in perspective, it took a team of researchers 900 core-years of computing time to factorise the 240-digit (795-bit) modulus n that is used in RSA-240 encryption.

These researchers estimated that increasing the number of bits to 1024 would make their factorisation program take 500 times as long to crack the encryption.

However, this does not mean there doesn't exist a fast algorithm for factorising integers. We just haven't found one yet (or if people have, they are keeping the algorithm to themselves). Most computer scientists believe that this is impossible though. It is conjectured that integer factorisation is a difficult problem—that is, there is no algorithm to factorise a b -bit integer in time polynomial in b . Until this conjecture is proved, all of the data we rely on encryption to protect could theoretically be in jeopardy.

III. GRAPH THEORY

*Quamobrem, cum nuper problematis cuiusdam mentio esset facta,
quod quidem ad geometriam pertinere videbatur,
at ita erat comparatum, ut neque determinationem quantitatum requieret,
neque solutionem calculi quantitatum ope admitteret,
id ad geometriam situs referre haud dubitavi:
praesertim quod in eius solutione solus situs in considerationem veniat,
calculus vero nullius prorsus sit usus.*

— LEONHARD EULER “*Solutio problematis ad geometriam situs pertinentis*” (1741)

12. Definitions and basic notions

A *graph* is a pair $G = (V, E)$ where V is a nonempty set and

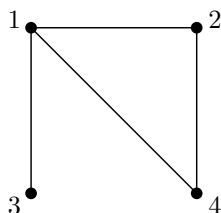
$$E \subseteq \{\{u, v\} : u, v \in V, u \neq v\}.$$

The graph G is said to be *finite* if and only if V is finite. The elements of V are called *vertices* and the elements of E are called edges. If $e = \{u, v\} \in E$, we say that u and v are *adjacent*, and e is *incident on* u and v . For brevity, we often write uv instead of $\{u, v\}$. Note that an edge is defined to be a 2-element set, not an ordered pair, so $uv = vu$.

Let's start with a small example. Let $V = \{1, 2, 3, 4\}$ and

$$E = \{12, 13, 14, 24\}.$$

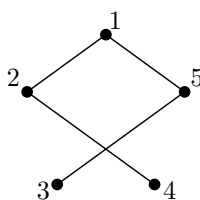
We can draw the graph $G = (V, E)$ as follows:



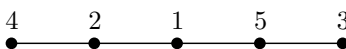
Instead of listing out all the vertices, we can define graphs by setting conditions for when a pair of vertices should be connected by an edge. For instance, let $G = (V, E)$ be the graph with $V = \mathbf{Z}/5\mathbf{Z}$ and, for $a, b \in V$ with $a \neq b$, we put $ab \in E$ if and only if

$$a + b \equiv 1 \pmod{5} \quad \text{or} \quad a + b \equiv 3 \pmod{5}.$$

Here is a drawing of G :



Turns out this graph was simpler than it might have seemed at first glance. Here is an alternate drawing of G :



Note that when we refer to G , we are talking about a set of vertices and edges, not any particular drawing! There are infinitely many valid drawings of any graph.

Suppose $G = (V, E)$ has a finite number of vertices. Then we can put the elements of V in bijection with $\{1, \dots, n\}$. Without loss of generality, we can simply assume that $\{1, \dots, n\}$ is the vertex set of G . The *adjacency matrix* of G is then defined to be the symmetric $n \times n$ matrix obtained by setting

$$A_{ij} = \begin{cases} 1, & \text{if } ij \in E, \\ 0, & \text{otherwise.} \end{cases}$$

For instance, the adjacency matrix of the graph on $\mathbf{Z}/5\mathbf{Z}$ defined earlier is

$$A = \begin{pmatrix} 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \end{pmatrix}.$$

We now define certain important families of graphs. For any integer $n \geq 1$, we define the *complete graph* K_n to be the graph on n vertices with every possible edge present. Small examples are illustrated in Fig. 2.

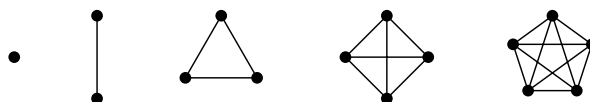


Fig. 2. Complete graphs K_n for small n .

31.X For an integer $n \geq 1$, we define the *Hamming cube* or *hypercube* to be the graph Q_n whose vertex set is $\{0, 1\}^n$, the set of all binary strings of length n , and whose edge set is the exactly the set of pairs of strings that differ by exactly one bit. Small examples are illustrated in Fig. 3. (As an exercise, try to figure out which strings correspond to which vertices!)

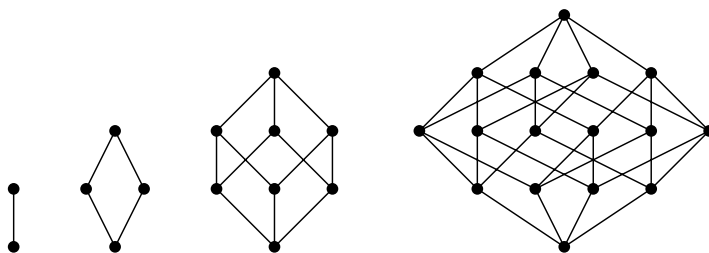


Fig. 3. Hypercubes Q_n for small n .

Degrees and k -regularity. The *neighbours* of a vertex v are all $u \in V$ such that $uv \in E$. The *degree* of a vertex v , denoted by $\deg(v)$, is the number of neighbours that v has. A graph is said to be *k -regular* for some $k \in \mathbf{N}$ if every $v \in V$ has degree k . For example, let C_n denote the graph on $V = \mathbf{Z}/n\mathbf{Z}$ where we connect $a \in V$ and $b \in V$ if and only if $a - b \equiv 1 \pmod{n}$ or $a - b \equiv -1 \pmod{n}$. Every vertex a has exactly two neighbours, namely, the congruence classes of $a - 1$ and $a + 1$ modulo n . So C_n is 2-regular for all n . In fact, complete graphs and hypercubes are also regular, for in K_n , each vertex is connected to each of the $n - 1$ other vertices, and in Q_n , each string is connected to the n other strings obtained by flipping each of its bits.

The following theorem relates vertex degrees to the number of edges.

Theorem 60. *Let $G = (V, E)$ be a finite graph. Then*

$$\sum_{v \in V} \deg(v) = 2|E|.$$

Proof. Consider each edge $e = uv \in E$. It contributes $+2$ to the right-hand side. On the left-hand side, it contributes $+1$ to the degree $\deg(u)$ of u , and $+1$ to the degree $\deg(v)$ of v . ■

The proof of this theorem is an excellent example of a proof by double counting, where we have proved an identity by showing that both sides are different ways of counting the same thing. This theorem has the following corollary.

Corollary 61 (*Handshaking lemma*). *In every finite simple graph, the number of vertices having odd degree is even.*

Proof. We partition $V = V_o \cup V_e$, where V_o comprises all vertices of odd degree and V_e comprises all vertices of even degree. By the previous theorem, we have

$$2|E| = \sum_{v \in V} \deg(v) = \sum_{v \in V_o} \deg(v) + \sum_{v \in V_e} \deg(v).$$

Taking this whole identity modulo 2, we have

$$0 \equiv \sum_{v \in V_o} 1 + \sum_{v \in V_e} 0 \pmod{2}.$$

Hence

$$|V_o| = \sum_{v \in V_o} 1 \equiv 0 \pmod{2},$$

which is what we wanted. ■

This proposition is often called the handshaking lemma (sometimes this is also used to refer to the preceding theorem), since it implies that at any gathering, the number of people who shake hands with an odd number of people is even. From Theorem 60 we also able to derive a corollary that counts the number of edges in k -regular graphs.

Corollary 62. *Let $G = (V, E)$ be k -regular. Then*

$$|E| = \frac{k|V|}{2}.$$

Proof. By Theorem 60, we have

$$2|E| = \sum_{v \in V} \deg(v),$$

but $\deg(v) = k$ for all $v \in V$ by k -regularity. Hence

$$2|E| = \sum_{v \in V} k = k|V|,$$

and the result follows upon rearranging. ■

Walks, paths, and cycles. A *walk* in $G = (V, E)$ is a sequence σ of vertices

$$\sigma = (v_0, v_1, \dots, v_n)$$

such that $v_i v_{i+1} \in E$ for all $0 \leq i \leq n-1$. The *endpoints* of σ are v_0 and v_n , and the *length* of σ , denoted $|\sigma|$, is the number of edges traversed, namely, n . The walk is said to be *closed* if $v_0 = v_n$ and *open* otherwise. A walk is called a *path* if no vertex is repeated.

Theorem 63. *Let $G = (V, E)$ be a graph. If u and v are vertices such that there exists a walk from u to v , then there exists a path from u to v .*

Proof. We perform a minimality argument. Let $\sigma = (v_0, v_1, \dots, v_n)$ be a walk from u to v (so $u = v_0$ and $v = v_n$) of shortest length. We claim that σ is in fact a path. Indeed, suppose for a contradiction that it is not a path; then there is some repeated vertex, so there exist $i, j \in \{0, 1, \dots, n\}$ such that $i < j$ and $v_i = v_j$. Hence there is no need to visit any of the vertices between v_i and v_j in σ , since $v_i = v_j$ is connected to v_{j+1} . Concretely, consider the walk

$$\sigma' = (v_0, v_1, \dots, v_i, v_{j+1}, \dots, v_n).$$

Note that $|\sigma'| = |\sigma| - (j - i)$, and $j - i > 0$, so σ' is a shorter walk from u to v . But this contradicts our choice of σ as a walk from u to v of shortest length. We conclude that σ is a path. ■

A *cycle* is a walk $\sigma = (v_0, v_1, \dots, v_n)$ of length at least 3 with $v_0 = v_n$ and no vertex repeated except $v_0 = v_n$.

Proposition 64. *Let $G = (V, E)$. If G contains a closed walk of odd length, then it contains a cycle of odd length.*

Proof. The idea is similar to the proof of the previous theorem. We perform a minimality argument. Let $\sigma = (v_0, v_1, \dots, v_n)$ be an odd-length closed walk in G , and choose this walk to have minimal odd length (i.e, any shorter walk has even length). We shall prove that σ is a cycle.

For a contradiction, suppose that σ is not a cycle, so that there exist indices i and j with $i < j$ such that $v_i = v_j$. Consider the two closed walks

$$\sigma_1 = (v_0, v_1, \dots, v_i, v_{j+1}, \dots, v_n)$$

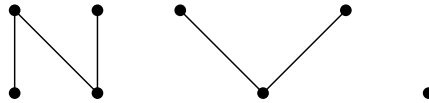
and

$$\sigma_2 = (v_i, v_{i+1}, \dots, v_j).$$

Both are shorter than σ , so by the minimality of σ , $|\sigma_1|$ and $|\sigma_2|$ are even. But this implies that $|\sigma| = |\sigma_1| + |\sigma_2|$ is even. This contradicts the hypothesis that $|\sigma|$ is odd.

We conclude that σ is a cycle. ■

05.XI **Connectedness.** We say a graph $G = (V, E)$ is *connected* if for all $u, v \in V$ there exists a walk (and hence a path) from u to v . All of the graphs we have drawn so far have been connected. Here is a graph that is disconnected:



This picture is meant to depict one graph on 8 vertices, not three different graphs.

The small examples of K_n and Q_n in Figs. 2 and 3 were all connected, so it is natural to assume that K_n and Q_n are connected for all n . This is indeed the case.

Proposition 65. *For all $n \geq 1$, the graphs K_n and Q_n are connected.*

Proof. Any vertices u and v in the vertex set of K_n belong to the edge uv , so we have the path (u, v) of length 1 between u and v . This shows that K_n is connected.

Now let u and v be any two vertices in K_n . Suppose there are m bits that differ between u and v . Then we can flip them one by one to change u into v . This represents a path of length m from u to v , since there is an edge in Q_n between any two strings in Q_n that differ at exactly one bit. ■

Let $G = (V, E)$ and let $u, v \in V$. We say that v is *reachable* from u if there exists a walk from u to v . If we define a relation \sim on V by setting $u \sim v$ if and only if v is reachable from u , then \sim is an equivalence relation (prove this as an exercise!) and the equivalence classes are called the *connected components* of G .

Here is an example. Let $G = (\mathbf{Z}, E)$, where for $i, j \in \mathbf{Z}$ with $i < j$, we have $ij \in E$ if and only if $j - i \in \{9, 15\}$. That is, every edge either hops forwards or backwards in the number line in increments of 9 or 15. Starting at $n \in \mathbf{Z}$, we can reach any vertex m that is of the form

$$m = n + 15s + 9t$$

for some $s, t \in \mathbf{Z}$. By Proposition 37 (this was the proposition about frogs and lilypads, if that jogs your memory), the integers representable as $15s + 9t$ for some $s, t \in \mathbf{Z}$ are exactly the multiples of $\gcd(9, 15) = 3$. So in fact, from n one can reach any integer m of the form $n + 3l$ for some $l \in \mathbf{Z}$; that is, one can reach any integer m with $m \equiv n \pmod{3}$. So the three connected components in this graph are $[0]_3$, $[1]_3$, and $[2]_3$, the equivalence classes of integers modulo 3.

13. Triangles and bipartite graphs

A *subgraph* of $G = (V, E)$ is a graph $G' = (V', E')$ such that $V' \subseteq V$ and $E' \subseteq E$ where for all $e = uv \in E'$, we have $u, v \in V'$. For example, the *triangle graph* C_3 is a subgraph of K_4 (draw it!). In fact, any graph on n vertices or less is a subgraph of K_n , since K_n contains all possible edges among n vertices.

We now consider an extremal question concerning subgraphs. An extremal question is one of the form, *What is the extremal (maximum or minimum) number of objects we can have, subject to some restriction?* The extremal question we ask is, *What is the maximum number of edges that a graph on n vertices can have, assuming that it does not contain a triangle as a subgraph?* To answer this question, we need the following inequality, which is one of the most important inequalities in all of mathematics. It is named for A. Cauchy, who in 1821 proved the version we shall record below, and H. Schwarz, who proved a version for integrals in 1888.

Theorem 66 (*Cauchy–Schwarz inequality*). For all $u_1, \dots, u_n, v_1, \dots, v_n \in \mathbf{R}$, we have

$$\left(\sum_{i=1}^n u_i v_i \right)^2 \leq \left(\sum_{i=1}^n u_i^2 \right) \left(\sum_{i=1}^n v_i^2 \right).$$

In the proof below, you will need to recall from your high school days that for a quadratic polynomial $ax^2 + bx + c$, the *discriminant* is defined to be the quantity $b^2 - 4ac$. (It appears under the square-root sign in the quadratic formula.) If it is negative, the polynomial has no real roots, if it is positive, the polynomial has two real roots, and if it is zero, the polynomial has one real root.

Proof. If $u_i = 0$ for all $1 \leq i \leq n$, then both sides are 0 and the inequality is trivially true. So suppose that at least one of the u_i is nonzero. Then

$$p(x) = (u_1 x + v_1)^2 + (u_2 x + v_2)^2 + \cdots + (u_n x + v_n)^2$$

is a quadratic polynomial in the variable x . Rewriting $p(x)$ in summation notation, we have

$$\begin{aligned} p(x) &= \sum_{i=1}^n (u_i x + v_i)^2 \\ &= \sum_{i=1}^n (u_i^2 x^2 + 2u_i v_i x + v_i^2) \\ &= \left(\sum_{i=1}^n u_i^2 \right) x^2 + 2 \left(\sum_{i=1}^n u_i v_i \right) x + \left(\sum_{i=1}^n v_i^2 \right). \end{aligned}$$

Since it is defined as a sum of squares, $p(x) \geq 0$ for all $x \in \mathbf{R}$, so it has at most one real root. This means its discriminant is nonpositive; in other words,

$$4 \left(\sum_{i=1}^n u_i v_i \right)^2 - 4 \left(\sum_{i=1}^n u_i^2 \right) \left(\sum_{i=1}^n v_i^2 \right) \leq 0.$$

Dividing through by 4 and then rearranging gives exactly the inequality we sought. **■**

The Cauchy–Schwarz inequality can also be proved using some results you might know from linear algebra. Defining vectors $u, v \in \mathbf{R}^n$ by $u = (u_1, \dots, u_n)$ and $v = (v_1, \dots, v_n)$, the left-hand side of the Cauchy–Schwarz inequality is $(u \cdot v)^2$ and the right-hand side is $\|u\|^2 \|v\|^2$. But, letting θ is the angle between u and v , then $u \cdot v = \cos \theta \|u\| \|v\|$, and applying the fact that $|\cos \theta| \leq 1$, we easily deduce the desired inequality.

We are now able to derive the maximum number of edges in a triangle-free graph on n vertices. The following theorem was proved by W. Mantel in 1907.

Theorem 67. (*Mantel’s theorem*). *Let $G = (V, E)$ be a graph not containing a triangle as a subgraph. Then*

$$|E| \leq \left\lfloor \frac{|V|^2}{4} \right\rfloor.$$

Proof. Consider the sum

$$\sum_{uv \in E} (\deg(u) + \deg(v)).$$

The term $\deg(u)$ appears in the sum exactly once for every edge incident on u ; that is, it appears $\deg(u)$ times. This is true for all $u \in V$, so we conclude that

$$\sum_{uv \in E} (\deg(u) + \deg(v)) = \sum_{u \in V} \deg(u)^2.$$

On the other hand, since G contains no triangle, for every pair of vertices u and v , the set of neighbours of u is disjoint from the set of neighbours of v . So $\deg(u) + \deg(v) \leq |V|$, and we have

$$\sum_{u \in V} \deg(u)^2 = \sum_{uv \in E} (\deg(u) + \deg(v)) \leq |V| \cdot |E|.$$

By Theorem 60 and the Cauchy–Schwarz inequality in that order, we have

$$\begin{aligned} (2|E|)^2 &= \left(\sum_{u \in V} \deg(u) \right)^2 \\ &= \left(\sum_{u \in V} \deg(u) \cdot 1 \right)^2 \\ &\leq \left(\sum_{u \in V} \deg(u)^2 \right) \left(\sum_{u \in V} 1^2 \right) \\ &= |V| \left(\sum_{u \in V} \deg(u)^2 \right). \end{aligned}$$

Hence

$$4|E|^2 \leq |V| \left(\sum_{u \in V} \deg(u)^2 \right) \leq |V|^2 \cdot |E|.$$

Rearranging this gives us

$$|E| \leq \frac{|V|^2}{4},$$

and we can take the floor function on the right-hand side, since $|E|$ must be an integer. ■

So if $G = (V, E)$ has $|E| \geq \lfloor |V|^2/4 \rfloor$, there must be a triangle subgraph in G . Now let $|V|$ be even, so that $|V|^2/4$ is an integer. Does there exist a graph G with exactly $|V|^2/4$ edges such that G does not contain a triangle as a subgraph? This brings us to the definition of a bipartite graph.

Bipartite graphs. A graph $G = (V, E)$ is called *bipartite* if there is a partition $V = A \cup B$ of the vertex set (so $A \cap B = \emptyset$) called the *bipartition* such that each edge has one endpoint in A and one endpoint in B . For example, hypercubes are bipartite.

Proposition 68. *For all $n \geq 1$, the hypercube Q_n is bipartite.*

Proof. Let $Q_n = (V, E)$. Every element $s \in V$ corresponds to a string $S = (s_1, \dots, s_n)$ of bits, where $s_i \in \{0, 1\}$ for all $1 \leq i \leq n$. Define

$$V_0 = \{s \in V : s_1 + \dots + s_n \equiv 0 \pmod{2}\}$$

and

$$V_1 = \{s \in V : s_1 + \dots + s_n \equiv 1 \pmod{2}\}.$$

It is clear that $V_0 \cup V_1 = V$ and $V_0 \cap V_1 = \emptyset$, so this is a bipartition of the vertex set, and for every $e = s_1 s_2 \in E$, the strings s_1 and s_2 differ in exactly one bit, so if $s_1 \in V_0$, then $s_2 \in V_1$, and vice versa. ■

The *complete bipartite graph* $K_{m,n}$ is a graph with $V = V_m \cup V_n$, where $|V_m| = m$ and $|V_n| = n$, and E is the set $\{uv : u \in V_m, v \in V_n\}$ of all possible edges between V_m and V_n . Fig. 4 illustrates two small examples of complete bipartite graphs.

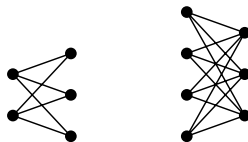


Fig. 4. The complete bipartite graphs $K_{2,3}$ and $K_{4,3}$.

When n is even, the graph $K_{n/2, n/2}$ has n vertices and exactly $n^2/4$ edges. When n is odd, the graph $K_{(n+1)/2, (n-1)/2}$ has n vertices and

$$\frac{n+1}{2} \cdot \frac{n-1}{2} = \frac{n^2-1}{4} = \frac{n^2}{4} - \frac{1}{4} = \left\lfloor \frac{n^2}{4} \right\rfloor$$

edges. (The last equality here follows from the fact that any odd n is congruent to either 1 or 3 modulo 4, which means its square is congruent to 1 modulo 4.)

The calculation above shows that $K_{\lceil n/2 \rceil, \lfloor n/2 \rfloor}$ always has n vertices and $\lfloor n^2/4 \rfloor$ edges. It is the largest possible number of edges such that Mantel's theorem does not apply. In other words, from Mantel's theorem alone, we do not know whether there must be a triangle in $K_{\lceil n/2 \rceil, \lfloor n/2 \rfloor}$, but adding a single edge to $K_{\lceil n/2 \rceil, \lfloor n/2 \rfloor}$ results in a graph that must have a triangle by Mantel's theorem.

We will actually be able to prove that $K_{\lceil n/2 \rceil, \lfloor n/2 \rfloor}$ does not contain any triangles. This means that the lower bound of $\lfloor n^2/4 \rfloor$ in Mantel's theorem is the best possible one.

First we need a lemma.

07.XI **Lemma 69.** *A graph $G = (V, E)$ is bipartite if and only if all of its connected components are bipartite.*

Proof. Let $G_1 = (V_1, E_1), G_2 = (V_2, E_2), \dots, G_n = (V_n, E_n)$ be the connected components of G .

If each connected component is bipartite, say, $V_i = A_i \cup B_i$ is a bipartition of V_i for all $1 \leq i \leq n$, then $A = \bigcup_{i=1}^n A_i$ and $B = \bigcup_{i=1}^n B_i$ is a bipartition of V , and every $e \in E$ has one endpoint in A and one endpoint in B .

Now suppose that there is some component G_i that is not bipartite. Now let $V = A \cup B$ be any partition of V into disjoint nonempty sets. Then $A_i = A \cap V_i$ and $B_i = B \cap V_i$ form a partition of V_i into two nonempty sets. But G_i is not bipartite, so there is some $e \in E$ that either has both endpoints in A_i or both

endpoints in B_i . This means that $A \cup B$ is not a valid bipartition, and since A and B were arbitrary, we conclude that G is not bipartite either. ■

As a matter of fact, we shall prove something much stronger than just the fact that $K_{\lceil n/2 \rceil, \lfloor n/2 \rfloor}$ does not contain any triangles. The following theorem characterises the family of bipartite graphs as *exactly* those graphs not containing odd-length cycles. Its proof shall require the following metric imposed on a graph G . For all $u, v \in V$ define the *distance* $\text{dist}(u, v)$ between u and v to be the length of the shortest path from u to v . If no such path between u and v exists, we set $\text{dist}(u, v) = \infty$.

Theorem 70. *A graph is bipartite if and only if it does not contain any cycles of odd length.*

Proof. First, assume that $G = (V, E)$ is bipartite; let $V = A \cup B$ be the bipartition of the vertex set. Let σ be a cycle in G . Each edge changes side (either A to B or B to A), so in order for the starting and ending vertices of this cycle to be the same, $|\sigma|$ must be even.

Now assume that $G = (V, E)$ has no odd-length cycles. To show that G is bipartite, it suffices to show that each of its connected components is bipartite (by the previous lemma), so without loss of generality we may assume that G is connected. This means that $\text{dist}(u, v) < \infty$ for all $u, v \in V$.

Select one vertex $h \in V$ and set

$$V_0 = \{v \in V : \text{dist}(h, v) \equiv 0 \pmod{2}\}$$

and

$$V_1 = \{v \in V : \text{dist}(h, v) \equiv 1 \pmod{2}\}.$$

Let $e = uv \in E$. Consider a closed walk σ that follows a shortest path from u to h , then a shortest path from h to v , and finally the edge e going from v to u . We have

$$|\sigma| = \text{dist}(u, h) + \text{dist}(h, v) + 1.$$

But since G contains no odd-length cycles, by (the contrapositive of) Proposition 64, any closed walk in G must have even length. This means that $\text{dist}(u, h)$ and $\text{dist}(h, v)$ are not the same modulo 2; that is, either $u \in V_0$ and $v \in V_1$ or vice versa. ■

14. Trees

A graph G is called a *forest* if it has no cycles. A connected forest is called a *tree*. A vertex of degree 1 in a forest is called a *leaf*. Fig. 5 contains an example of a tree.

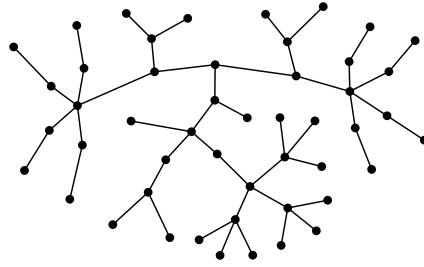


Fig. 5. A tree with 25 leaves (count them!).

Our first proposition concerning trees is that every tree containing at least one edge has at least two leaves.

Proposition 71. *Let F be a forest with at least one edge (and hence at least two vertices). There are at least two leaves in F .*

Proof. Let $\sigma = (v_0, v_1, \dots, v_n)$ be a path of maximum length in F . Since there is at least one edge in F , we have $|\sigma| \geq 1$. The claim is that $\deg(v_0) = \deg(v_n) = 1$. Indeed, for a contradiction, assume that one of these two endpoints has degree at least two. Without loss of generality, suppose it is v_0 (we can reverse the path if needed to make this the case). So there exists $u \in V$ such that uv_0 is an edge. If $u = v_i$ for some $1 \leq i \leq n$, then there is a cycle $(v_i, v_0, v_1, \dots, v_i)$, contradicting the fact that F is a forest. So for all $1 \leq i \leq n$, $u \neq v_i$. This means that (u, v_0, \dots, v_n) is a strictly longer path than σ , contradicting the maximality of σ . ■

The following gives a characterisation of trees.

Proposition 72. *A graph $G = (V, E)$ is a tree if and only if for all $u, v \in V$, there exists a unique path from u to v in G .*

Proof. First we assume that G is a tree. Let $u, v \in V$. Since G is a tree, it is connected, so there must exist some path σ from u to v in G . Now we prove that this path is unique. For a contradiction, let σ' be another, different, path from u to v in G . Both σ and σ' start at u , but eventually they are different. Let v_1 be the last vertex they have in common (this vertex could be u). Eventually, σ and σ' must meet again, since they both end at v . Let v_2 be the first time they meet again, after v_1 (we could have $v_2 = v$). Now let σ be the path obtained by following σ from v_1 to v_2 , then following σ' backwards to v_1 . This is a cycle, by choice of v_1 and v_2 . But this gives a contradiction, since there are no cycles in G .

Now suppose that for all $u, v \in V$, there exists a unique path from u to v in G . The fact that there exists a path at all between every two vertices immediately implies that G is connected. It remains to prove that there is no cycle in G . But this is clear, since if σ is a cycle in G , then picking any two distinct vertices

in σ , we have two different paths between those two vertices (one obtained by following σ clockwise, and one obtained by following it anticlockwise). ■

Another characterisation of trees is that they are exactly the connected graphs on n vertices with $n - 1$ edges. To prove this, we first need a lemma.

12.XI **Lemma 73** (*Detour lemma*). *Let G be a connected graph and let σ be a cycle in G . If G' is the graph obtained by deleting one edge of σ , then G' is also connected.*

Proof. Let $e = v_1v_2$ be the deleted edge. Let u, v be any two vertices in G . Since G is connected, there exists a walk τ from u to v in G . We want to show that there is a path from u to v in G' . If τ does not use the edge e , then it is also a valid path in G' and we are done. If it does use e , then to get from v_1 to v_2 , we detour around σ in the other direction, then continue from v_2 to v , obtaining a walk in G' from u to w . ■

Theorem 74. *A graph $G = (V, E)$ is a tree if and only if $|E| = |V| - 1$ and G is connected.*

Proof. We prove the forward implication by induction on $n = |V|$, the number of vertices. Concretely, the statement we shall prove is, *For all $n \geq 1$, for all $G = (V, E)$ with $|V| = n$, if G is a tree then $|E| = |V| - 1$ and G is connected.* The base case is $n = 1$. Then the only possible graph G is simply \bullet , a single vertex. This is a tree since it is connected and contains no cycles. We have $|E| = 0 = 1 - 1 = |V| - 1$. For the inductive step, let $n \geq 1$ and assume the statement holds for n . Let $G = (V, E)$ be a graph with $|V| = n + 1$, and assume that G is a tree. Since $n \geq 1$, $|V| \geq 2$, so there are at least two vertices, and connectedness of G means that there is at least one edge in G . Hence G is a forest with at least one edge, and by Proposition 71, there are at least two leaves in G . Pick one of these leaves and call it u . Form a new graph $G' = (V', E')$ by removing u and the single edge incident on it; we have $|V'| = n$ and $|E'| = |E| - 1$. The graph G' is a tree, because for each $v, w \in V'$, there is a unique path from v to w in G , and this path does not pass through u (since u is a leaf). Hence

$$|E| = |E'| + 1 = |V'| - 1 + 1 = |V| - 1,$$

where in the second equality we have used the induction hypothesis.

Now assume that $|E| = |V| - 1$ and G is connected. We want to show that G has no cycles. For a contradiction, suppose it does; call it σ . Remove an edge from σ . By the detour lemma, this resulting graph is still connected. If this graph still has a cycle, remove an edge from it again, and repeat this process until no cycles remain. This graph G' still has V as its vertex set, but it has a new edge set E' with $|E'| < |E|$. But G' is now a tree, since it contains no cycles but is still connected, so by the previous paragraph, we have $|E'| = |V| - 1 = |E|$. The contradiction completes the proof. ■

We are now able to answer the following extremal question: *How many edges can a graph on n vertices have if it does not contain a cycle (i.e., if it is a forest)?*

Corollary 75. Let $F = (V, E)$ be a forest. Then $|E| \leq |V| - 1$.

Proof. Let $k \geq 1$ denote the number of connected components in F , and number the connected components C_1, C_2, \dots, C_k . Join C_1 to C_2 by an edge, C_2 to C_3 by an edge, and so on. This creates a connected graph $G = (V, E')$ with $|E'| = |E| + k - 1$, and G must be a tree, since adding these edges does not introduce a cycle. By the previous theorem, $|E'| = |V| - 1$, so $|E| \leq |V| - 1$. ■

15. Eulerian trails and circuits

Recall that a walk that does not repeat any vertex is called a path. A similar definition is that of a *trail*: this is a walk that does not repeat any *edge*. A walk is called an *Eulerian trail* if it uses every edge of G exactly once, and an *Eulerian circuit* if it is an Eulerian trail and it is closed (its endpoints are the same vertex).

The term “Eulerian” refers to L. Euler, who in 1735 answered the following question: *Is it possible to cross every bridge in Königsberg (a city in Prussia), without having to repeat any bridges?* The situation is illustrated in Fig. 6.

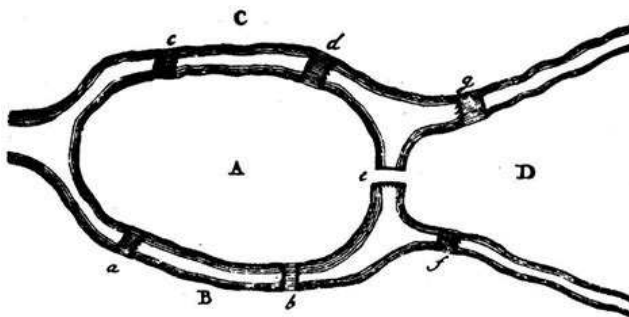


Fig. 6. The seven bridges of Königsberg in 1735.

To represent this map as a graph $G = (V, E)$, we will have to allow multiple edges between vertices. There are a couple ways to formalise this set-theoretically, but we will dispense with the details here and simply think of E as a special type of set that allows for copies of the same edge to be included multiple times. The graph G representing Königsberg is depicted in Fig. 7.

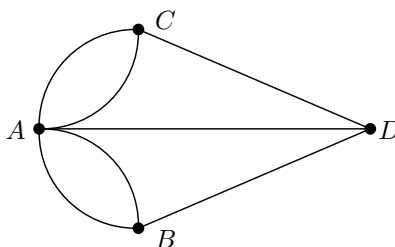


Fig. 7. The graph representing the seven bridges of Königsberg.

What Euler actually proves in his paper (which did not appear in print until 1741) is that it is impossible to cross every bridge in Königsberg exactly once; i.e., G does not contain an Eulerian trail. First we characterise the graphs that contain Eulerian circuits, as we'll need this to later characterise graphs with Eulerian trails.

Theorem 76. *Let $G = (V, E)$ be a connected graph with a finite number of vertices and edges, where multiple edges between the same two vertices are allowed. Then G has an Eulerian circuit if and only if all vertices have even degree.*

Proof. Suppose G has an Eulerian circuit $\sigma = (v_0, v_1, \dots, v_n)$ where $v_0 = v_n$. For each $v_i \neq v_0$, each visit uses two edges, one edge to enter v_i and one to exit it, so $\deg(v_i)$ is even for all $0 \leq i < n$. For $v_0 = v_n$, there is one edge that leaves it at the beginning, two edges for each visit during the circuit, and then one returning to it at the end, so $\deg(v_0)$ is even as well.

We prove the reverse implication by induction on $m = |E|$. That is, we prove that for all integers $m \geq 0$, if $G = (V, E)$ with $|E| = m$ and all vertices in G have even degree, then G has an Eulerian circuit.

For the base case $m = 0$, the stipulation that G be connected forces $G = \bullet$, and this graph has an Eulerian circuit (the trivial walk).

Now let $m \geq 0$ and suppose the statement is true for all integers less than or equal to m . Let $G = (V, E)$ with $|E| = m + 1$. Since there is at least one edge in G now, G cannot be a tree. This is because every tree with at least one edge contains at least two leaves, and leaves have (odd) degree 1. Let σ be a cycle in G and form $G' = (V, E')$ by removing every edge in σ . This graph may no longer be connected, but since we have removed a cycle, every vertex in G' either has the same degree as in G (if the vertex was not in σ), or degree two less than in G (if the vertex was in σ). So all vertices in G' have even degree. Let k be the number of connected components in G' ; number them H_1, \dots, H_k . Each has at most m edges, is connected, and all vertices have even degree, so each H_i contains an Eulerian circuit, by induction. We build an Eulerian circuit in G as follows.

- i) Start at any vertex of σ .
- ii) Follow σ .
- iii) Each time we reach a vertex of a connected component H_i of G' for the first time, follow its Eulerian circuit.
- iv) Once we return to the beginning of this circuit, continue along σ .
- v) Repeat this process until we return to the vertex of σ that we started at.

This gives an Eulerian circuit of G . ■

Now we characterise Eulerian trails.

Theorem 77. *Let $G = (V, E)$ be a graph (with multiple edges allowed). Then G contains an Eulerian trail that isn't an Eulerian circuit if and only if exactly two vertices of G have odd degree.*

Proof. Suppose that G contains an Eulerian trail $\sigma = (v_0, v_1, \dots, v_n)$ in which $v_0 \neq v_n$. By a similar reasoning to the first paragraph of the previous proof, every vertex in the walk has even degree, except both v_0 and v_n must have odd degree, since starting at v_0 we must leave for the first time via one edge, and at the end we enter v_n one last time without exiting.

Now assume that G contains exactly two odd-degree vertices. By the previous theorem, there does not exist an Eulerian circuit in G . Let u and v be the two odd-degree vertices. Add a new edge $e = uv$, even if there was already an edge between u and v before. We get G' in which every vertex has even degree. So G' must have an Eulerian circuit, call it σ . It contains the edge e , but we can remove e to get an Eulerian trail in the original graph G . ■

Since the graph representing the Königsberg bridge problem has four vertices, all of odd degree, the two theorems we just proved tell us that it contains neither an Eulerian trail nor an Eulerian circuit. Hence it is impossible to cross every bridge in Königsberg exactly once.

16. Planar graphs

- 14.XI A graph G is called *planar* if we can draw G in the plane without edges crossing. If we can, such a drawing is called a *planar embedding* of G . The graph K_4 is planar; its representation in Fig. 2 does not make this evident, but a planar embedding is given in Fig. 8. Likewise, the hypercube Q_3 was drawn with edge crossings in Fig. 3, but it is also planar.

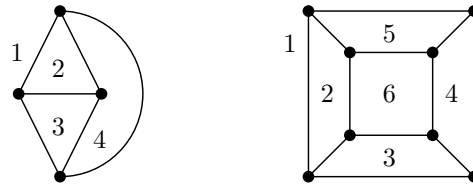


Fig. 8. Planar embeddings of K_4 and Q_3 , with numbered faces

Planarity is a property that passes down to subgraphs, as the following proposition shows.

Proposition 78. *If G is planar then any subgraph H of G is also planar.*

Proof. If G is drawn in the plane without edges crossing, and we erase some edges and vertices to create a drawing of H , we cannot introduce any edge crossings along the way. ■

A planar embedding of a graph G determines regions, or *faces* bounded by the edges, including one region “outside” the graph. We define the *Euler characteristic* of a planar embedding of $G = (V, E)$ to be the quantity

$$\chi = |V| - |E| + f,$$

where f is the number of faces in the planar embedding. For the examples in Fig. 8, we have

$$\chi = 4 - 6 + 4 = 2$$

in the case of K_4 , and

$$\chi = 8 - 12 + 6 = 2$$

in the case of Q_3 . This might seem like a funny coincidence. The following theorem shows that it is not.

Theorem 79 (*Euler's formula*). *Let $G = (V, E)$ be a connected planar graph, let f be the number of faces in a planar embedding of G . Then the Euler characteristic χ of G equals 2.*

The proof of this formula requires a theorem whose statement seems so ridiculously intuitive that it seems silly that it even needs to be stated as a theorem. Despite its innocuous appearance, is surprisingly difficult to prove, so much so that its original 1887 proof by C. Jordan was thought by many mathematicians to be so lacking in detail as to be wrong, and the first widely accepted proof did not appear until 1905 (due to O. Veblen). More recently some mathematicians have analysed Jordan's original proof and found that it was more or less correct to begin with, just not presented in a satisfactory way. In any case, here is the theorem.

Theorem 80 (*Jordan curve theorem*). *Every closed curve in the plane \mathbf{R}^2 that does not intersect itself divides the plane into two distinct regions.* ■

The two regions that the closed curve creates are simply the inside of the curve and the outside of the curve. Using this theorem, we can now prove Euler's formula.

Proof of Theorem 79. We proceed by induction on $m = |E|$. In the case that $m = 0$, G can only be \bullet , in which case we have $\chi = 1 - 0 + 1 = 2$.

Now let $m \geq 0$ and assume that the theorem holds for all natural numbers at most m . Let $G = (V, E)$ be a connected planar graph with $|E| = m + 1$. Draw G in the plane and let f be the number of faces in the embedding. There are two cases, according to whether G is a tree or not. If G is a tree, then $f = 1$, since G has no cycles. But we also have $|E| = |V| - 1$, so

$$\chi = |V| - |E| + f = |V| - (|V| - 1) + 1 = 2.$$

If G is not a tree, then G has a cycle, call it σ . Form G' by removing exactly one edge e of σ ; so

$$|E'| = |E| - 1 = m + 1 - 1 = m.$$

By the detour lemma, G' is still connected, and it is planar, since removing an edge doesn't create any crossings. By the inductive hypothesis, the Euler characteristic χ' of G' satisfies

$$\chi' = |V| - |E'| + f' = 2,$$

where f' is the number of faces in the drawing of G . Returning to the drawing of G , the cycle σ is a closed curve without intersections in the plane, so it divides the plane into two distinct regions (each of which may have multiple faces within them). In particular, the removed edge e forms a border between two distinct faces, and in G' , these two faces combine into one face. Hence $f = f' + 1$. Putting everything together, we calculate that the Euler characteristic χ of G satisfies

$$\begin{aligned}\chi &= |V| - |E| + f \\ &= |V| - (|E'| + 1) + (f' + 1) \\ &= |V| - |E'| + f' \\ &= 2. \quad \blacksquare\end{aligned}$$

Euler's formula tells us that the number of faces in any drawing of G depends entirely on the number of vertices and edges that G has. We can use this information to give conditions for a graph to be non-planar.

Theorem 81. *Let $G = (V, E)$ be a connected planar graph with $|V| \geq 5$. Then*

$$|E| \leq 3|V| - 6.$$

Under the further assumption that G contains no triangles, we have the better bound

$$|E| \leq 2|V| - 4.$$

Proof. The proof is by double counting. Let R be the set of all regions into which G divides the plane, so that $|R| = f$, the number of faces. Consider the set

$$S = \{(e, r) \subseteq E \times R : \text{the edge } e \text{ touches the region } r\}.$$

We will count $|S|$ in two ways. First off, each edge $e \in E$ touches at most two regions, so $|S| \leq 2|E|$. On the other hand, every region is bounded by a cycle, and cycles have at least 3 edges, so each region $r \in R$ touches at least 3 edges. In other words, $|S| \geq 3f$. Chaining these two inequalities, we have

$$3f \leq |S| \leq 2|E|.$$

But by Euler's formula, we have $|V| - |E| + f = 2$, so $f = 2 - |V| + |E|$, and substituting this above, we get

$$3(2 - |V| + |E|) \leq 2|E|.$$

Distributing and rearranging terms yields the desired inequality $|E| \leq 3|V| - 6$.

If we have the further assumption that there are no triangles in G , then every region must touch at least 4 edges, allowing us to conclude the stronger bound $4f \leq 2|E|$. Then we proceed as above to get

$$4(2 - |V| + |E|) \leq 2|E|,$$

which can be manipulated to give $|E| \leq 2|V| - 4$. \blacksquare

The condition that $|V| \geq 5$ is a complete non-issue, since every graph with fewer than 5 vertices is a subgraph of K_4 , and we can apply Proposition 78, since we already know that K_4 is planar. As a corollary of the (contrapositive of the) previous theorem, we can give two important examples of nonplanar graphs.

Corollary 82. *The complete graph K_5 and the complete bipartite graph $K_{3,3}$ are nonplanar.*

Proof. If $(V, E) = K_5$ we have $|V| = 5$ so $3|V| - 6 = 9$, but $|E| = 5 \cdot 4/2 = 10$. So by the contrapositive of the previous theorem K_5 cannot be planar.

The graph $K_{3,3}$ does not contain any triangles, so if $K_{3,3}$ is planar, then the stronger bound $|V| \leq 2|E| - 4$ must hold. But in $(V, E) = K_{3,3}$, we have $|V| = 6$, so $2|E| - 4 = 8$, and we have $|E| = 9$. So $K_{3,3}$ is also nonplanar. ■

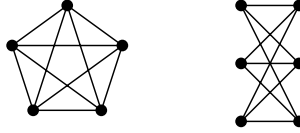


Fig. 9. The nonplanar graphs K_5 and $K_{3,3}$.

This corollary, combined with the contrapositive of Proposition 78, shows that K_n is nonplanar for all $n \geq 5$, and $K_{m,n}$ is non-planar whenever m and n are both at least 3. More generally, any graph that contains K_5 or $K_{3,3}$ as a subgraph is nonplanar.

However, subgraphs are not exactly the right notion to be considering when talking about planarity. To see why, consider the graph in Fig. 10.

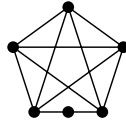


Fig. 10. The graph K_5 , with an extra vertex subdividing an edge.

It is nonplanar, since if it had a planar embedding, then by contracting one of the edges incident on the extra vertex and bending other edges accordingly, we would obtain a planar embedding of K_5 . On the other hand, K_5 is not a subgraph of this graph. We thus introduce the following notion. We say that a graph H is a *graph minor* of a graph G if H can be obtained from G by repeatedly either

- i) deleting an edge;
- ii) deleting a vertex; or
- iii) contracting an edge uv by removing it and merging u and v into a single vertex (and also combining any resulting multiple edges into a single edge).

Any number of the first two operations simply creates a subgraph. It is the contraction operation that produces interesting examples of graph minors. It is easy to see pictorially that contracting edges preserves planarity (we will not prove this rigorously), so *any graph minor of a planar graph is planar*. We can

use this to show, for instance, that the Petersen graph, depicted in Fig. 11, is nonplanar, since by contracting each of the five edges connecting the inner star pentagram to the outer pentagon, we obtain K_5 . (Note that Theorem 81 does not directly apply to the Petersen graph.)

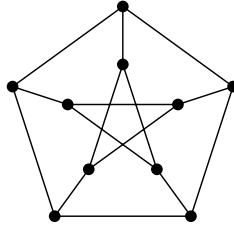


Fig. 11. The Petersen graph.

19.XI Not only are all graphs that have K_5 or $K_{3,3}$ as a minor nonplanar, it turns out that these are the only nonplanar graphs!

Theorem 83 (*Wagner's theorem*). A graph G is nonplanar if and only if either K_5 or $K_{3,3}$ is a graph minor of G . ■

This theorem was proven in 1937 by K. Wagner, but it is closely related to a theorem that was proved earlier (in 1930) by K. Kuratowski. It deals with the notion of a *subdivision* of a graph G : this is any graph obtained by repeatedly placing vertices on edges of G , subdividing an edge into two edges each time. (Fig. 10 was an example of a subdivision of K_5 .) We now state Kuratowski's theorem.

Theorem 84 (*Kuratowski's theorem*). A graph G is nonplanar if and only if it contains as a subgraph a subdivision of either K_5 or $K_{3,3}$. ■

It is easy to see that if H is a subdivision of G , then G is a minor of H , since we can reobtain G by contracting all the edges created by the subdivision operations. So the two theorems above are certainly very closely related, though it is not immediately obvious if one should imply the other. The truth is that they are equivalent, because it can be shown that any graph with either K_5 or $K_{3,3}$ as a minor also has a subgraph that is a subdivision of one of them.

IV. COMBINATORICS

*Quare merito suo utilissima censenda est Ars, Combinatoria dicta,
quae huic mentis nostrae defectui medetur,
docetque sic enumerare modos omnes posibles,
secundum quos res plures permisceri, transponi vel conjungi invicem possunt,
ut certi, simus, nos nullum eorum praetermisisse,
qui instituto nostro conducere valent.*

— JAKOB BERNOULLI, *Ars Conjectandi* (1713)

17. Counting techniques

Combinatorics is a branch of mathematics that, broadly speaking, deals with counting (or estimating) the cardinalities of finite sets. Let us now recollect the counting techniques that we know. If two sets A and B are disjoint, i.e., $A \cap B = \emptyset$, then the number of elements in the union $A \cup B$ is simply $|A| + |B|$. The cardinality of the power set $A \times B$ is $|A| \cdot |B|$, since for each of the $|A|$ elements $a \in A$, there are $|B|$ choices for $b \in B$ to pair with it. Lastly, if U is a universe and $A \subseteq U$, then \overline{A} is the set of all elements of the universe that don't belong to A , so $|\overline{A}| = |U| - |A|$.

Counting strings. Back in the proof of Proposition 5, we counted the number of binary strings of length n to be 2^n . More generally, if there are s different symbols and t different slots in which to put them, then the number of strings is s^t . This can be phrased in terms of the number of functions from a set T of size t to a set S of size s . Indeed, people often use the notation S^T to denote the set of all functions from T to S , and the notation should remind you that $|S^T| = |S|^{|T|}$. Identifying the number 2 with the set $\{0, 1\}$, the power set 2^X is in bijection with the set $\{0, 1\}^X$, so this notation for the set of all functions is consistent with our earlier notation for power set.

For example, using the symbol set $S = \{a, b, c, d\}$, the number of strings of length 5 is $|S|^5 = 4^5$. What about the number of strings of length *at most* 5, including one empty string of length 0? This is the number of strings of length 0, plus the number of strings of length 1, and so on up to and including length 5; that is, this equals

$$4^0 + 4^1 + 4^2 + 4^3 + 4^4 + 4^5.$$

To count strings of length 5 and containing at least one a , we let the universe be all strings of length 5 and let A be the set of all strings without any a , i.e., using only the letters b, c , and d . Then the number of strings containing at least one a is

$$|\overline{A}| = |U| - |A| = 4^5 - 3^5.$$

Next we count the number of strings of length 5, but we require that adjacent symbols be different. This one is best calculated by counting the number of choices we have at each slot. There are 4 choices for the first symbol, but after that we only have choices for each letter, since they cannot match the previous one. Hence the count is $4 \cdot 3^4$.

Lastly, to count the number of strings that must have at least one of the patterns aa , bb , cc , or dd , we simply note that this is the complement of what we counted in the previous paragraph, so the answer is $4^5 - 4 \cdot 3^4$.

The principle of inclusion and exclusion. As our next toy problem, suppose we want to count the number of strings of length 6, using the symbols a, b, c , and d , that start with a or end with aa (or both). The number of strings that start with a is 4^5 , since we have no choice for the first symbol, but 4 choices for each of the remaining five symbols; likewise, the number of strings that end

with aa is 4^4 . Adding these two up, we get $4^5 + 4^4$. There's a problem though. We have accidentally double-counted the strings that both start with a and end with aa . How many of these are there? Well, the first, fifth, and sixth slots all have to be a , but there are four choices for each of the remaining three slots, for a total of 4^3 . Subtracting this from our subtotal yields the correct answer: $4^5 + 4^4 - 4^3$.

This technique works with any two sets A and B , i.e., $|A \cup B| = |A| + |B| - |A \cap B|$. For three sets A , B , and C , we have

$$|A \cup B \cup C| = |A| + |B| + |C| - |A \cap B| - |A \cap C| - |B \cap C| + |A \cap B \cap C|$$

(draw a Venn diagram and convince yourself that this is correct). For a general number of sets, we have the following theorem.

Theorem 85 (*Principle of inclusion and exclusion*). *Let $n \geq 2$ be an integer and let A_1, \dots, A_n be finite sets. Then*

$$\left| \bigcup_{i=1}^n A_i \right| = \sum_{k=1}^n (-1)^{k+1} \left(\sum_{1 \leq i_1 < \dots < i_k \leq n} |A_{i_1} \cap \dots \cap A_{i_k}| \right).$$

This sum of sums might look intimidating, but breaking down what it says bit by bit, we see that it is simply a sum over all possible nonempty intersections of the A_i , with a sign according to how many of the A_i we are intersecting. We will prove this theorem soon, after we state and prove the binomial theorem. We first define, for nonnegative integers n and k , the *binomial coefficient* $\binom{n}{k}$ to be the number of ways of choosing k distinct elements out of a pool of n distinct elements. (The quantity $\binom{n}{k}$ is read “ n choose k .”) We have the formula

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

This is because we can pick k elements out of an n -element set by permuting all the n elements (there are $n!$ ways to do this), and then simply looking at the first k elements. But the order of the k elements we choose doesn't matter, so we arrive at the same elements no matter which of the $k!$ possible orders they might appear in. Same goes for the $(n-k)!$ orders the elements that we *didn't* pick can appear in.

The name “binomial coefficient” is justified by the following theorem, first stated and proved by Al-Karaji in around the tenth century.

Theorem 86 (*Binomial theorem*). *For all $x, y \in \mathbf{R}$ and positive integers n ,*

$$(x+y)^n = \sum_{k=0}^n \binom{n}{k} x^k y^{n-k}.$$

Proof. Expanding the left-hand side

$$(x+y)^n = (x+y)(x+y) \cdots (x+y),$$

we see that there will be 2^n different terms after repeatedly using the distributive law. Each term will be a product of x to some power and y to some power, where the powers add up to n . In other words, each term looks like $x^k y^{n-k}$ for some $0 \leq k \leq n$. The number of times that the term $x^k y^{n-k}$ appears in the expansion is the number of ways to choose k items out of a pool of n , namely, $\binom{n}{k}$. This gives exactly the right-hand side. ■

We are now able to prove the principle of inclusion and exclusion.

Proof of Theorem 85. Let $n \geq 2$ and let A_1, \dots, A_n be finite sets. Recall that the identity we want to prove is

$$\left| \bigcup_{i=1}^n A_i \right| = \sum_{k=1}^n (-1)^{k+1} \left(\sum_{1 \leq i_1 < \dots < i_k \leq n} |A_{i_1} \cap \dots \cap A_{i_k}| \right).$$

Let x be an element that belongs to the union $\bigcup_{i=1}^n A_i$. Clearly, this contributes 1 to the left-hand side, so we must show that it contributes exactly 1 to the right-hand side as well. We know that x is a member of at least one of the A_i ; let $s \geq 1$ be the number of sets A_i that contain x .

Each time x is a member of some A_i , it contributes +1 to the right-hand side. This happens $\binom{s}{1}$ times. Each time it is a member of A_i and A_j for $i < j$, x contributes -1 to the right hand side; there are $\binom{s}{2}$ of these terms. Generally, if x is a member of $A_1 \cap A_2 \cap \dots \cap A_r$, it contributes +1 if r is odd and -1 if r is even, and this term is repeated $\binom{s}{r}$ times. Hence x contributes

$$\binom{s}{1} - \binom{s}{2} + \binom{s}{3} - \dots + (-1)^{s+1} \binom{s}{s}$$

to the right-hand side.

But, by the binomial theorem,

$$0 = (1 - 1)^s = \binom{s}{0} - \binom{s}{1} + \binom{s}{2} - \dots + (-1)^s \binom{s}{s}.$$

We know that $\binom{s}{0} = s!/(0!s!) = 1$, so by rearranging we find that

$$\binom{s}{1} - \binom{s}{2} + \binom{s}{3} - \dots + (-1)^{s+1} \binom{s}{s} = \binom{s}{0} = 1,$$

which means that the contribution of x to the right-hand side is 1. We are done since x is arbitrary. ■

As an application of this principle, let $X = \{1, 2, \dots, 100\}$ and suppose we want to count the number of $n \in X$ with $\gcd(n, 30) = 1$. To accomplish this, it turns out to be easier to count the number of $n \in X$ that have $\gcd(n, 30) \geq 2$ (and then we must subtract this number from $|X| = 100$). For any positive integer r , let

$$A_r = \{n \in X : r \mid n\}.$$

Since $30 = 2 \cdot 3 \cdot 5$, an integer $n \in X$ has $\gcd(n, 30) \geq 2$ if and only if $n \in A_2$, $n \in A_3$, or $n \in A_5$. By the principle of inclusion and exclusion,

$$|A_2 \cup A_3 \cup A_5| = |A_2| + |A_3| + |A_5| - |A_2 \cap A_3| - |A_2 \cap A_5| - |A_3 \cap A_5| + |A_2 \cap A_3 \cap A_5|.$$

The number of even integers in X is $100/2 = 50$; similarly, the number of integers divisible by 3 in X is $\lfloor 100/3 \rfloor = 33$. In general, $|A_r| = \lfloor 100/r \rfloor$. Furthermore, since $\gcd(2, 3) = 1$, the intersection $A_2 \cap A_3$ is simply A_6 , and by analogous reasoning, we see that

$$\begin{aligned} |A_2 \cup A_3 \cup A_5| &= |A_2| + |A_3| + |A_5| - |A_6| - |A_{10}| - |A_{15}| + |A_{30}| \\ &= \left\lfloor \frac{100}{2} \right\rfloor + \left\lfloor \frac{100}{3} \right\rfloor + \left\lfloor \frac{100}{5} \right\rfloor - \left\lfloor \frac{100}{6} \right\rfloor - \left\lfloor \frac{100}{10} \right\rfloor - \left\lfloor \frac{100}{15} \right\rfloor + \left\lfloor \frac{100}{30} \right\rfloor \\ &= 50 + 33 + 20 - 16 - 10 - 6 + 3 \\ &= 74. \end{aligned}$$

We conclude that the number of $n \in X$ such that $\gcd(n, 30) = 1$ is $100 - 74 = 26$.

18. Permutations and combinations

Suppose we have a pool of 100 hockey players and we want to dress 20 of them for a game, giving them the numbers 1 through 20. How many ways are there to do this? Well, there are 100 choices for the player with jersey number 1, then after picking that player, there are 99 choices for the player with jersey number 2, and so on, so the total number is

$$100 \cdot 99 \cdot 98 \cdots 82 \cdot 81 = \frac{100!}{80!}.$$

21.XI In general, we define a *k-permutation* of a set X to be an ordered subset of X (i.e., the order of the selection matters) of size k . If $|X| = n$, the number of k -permutations of X is denoted $P(n, k)$ and equals

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

This number is also sometimes called the *falling factorial* and denoted $n^{\underline{k}}$. It is the same as the number of functions from $\{1, 2, \dots, k\}$ to X that are injective. (Pause a moment to think about it if you need to.) Recall that if $|S| = s$ and $|T| = t$, the total cardinality of S^T is s^t . Well, the number of *injective* functions from T to S is $P(s, t) = s^{\underline{t}}$.

An n -permutation of a set of size n is also just called a *permutation* of n . There are

$$P(n, n) = \frac{n!}{(n-n)!} = n!$$

of these. (Note that $0! = 1$, by definition.)

To illustrate various counting techniques associated with permutations, consider the following situation. We have ten books: three comic books, three books of poetry, and four novels. How many ways are there to order all of them on a shelf? The answer is the number of permutations on 10 objects, namely, $10! = 3628800$. How about just putting five of them on the shelf? Then the answer is $P(10, 5) = 10!/5! = 30240$. Lastly, what if we want to put all of them on the shelf, but group all the comic books together, all the poetry books together, and all the novels together. There are $3!$ ways to order the comic books amongst themselves, $3!$ for the poetry books as well, and then $4!$ for the novels. Then we must arrange the ordered bundles of books on the shelf; there are three bundles so there are $3!$ ways to do this. Thus the total number of ways is $3!3!4!3! = 5184$.

Returning to the hockey-team example, suppose we still want to pick a team of 20 out of 100 possible players, but now we don't number the players; that is, we only care about the sets of the players themselves, not the ordering amongst them. The number $P(100, 20)$ is $20!$ times too big, since there are $20!$ different orderings of the team we choose, and each of them is only counted once in our new scenario. So the answer is

$$\frac{P(100, 20)}{20!} = \frac{100!}{20!(100 - 20)!} = \binom{100}{20}.$$

We already had an intimate encounter with the quantity $\binom{n}{k}$ in the previous section. Defining a *k-combination* of a set X to be the number of ways of choosing k elements of X , the number of k -combinations of a set of size n is $\binom{n}{k}$.

Counting grid paths. An application of k -combinations is counting the number of paths in a grid. If you have ever had the pleasure or misfortune of visiting Edmonton, Alberta, you will know that most of the city's roads are organised in a numbered grid pattern. Avenues run east to west and streets run north to south. Suppose we have just enjoyed a cortado in the trendy café at the corner of 82nd Avenue and 104th Street, and now we want to get to the sandwich shop at the corner of 85th Avenue and 109th Street. To do so we will have to walk three blocks north and five blocks west. How many ways are there to do this?

Well, the journey is eight blocks in length, so each possible path can be expressed as a string of eight symbols, call them N for north and W for west. We also know there must be exactly three Ns and five Ws; for example, one possible path, pictured in Fig. 12, is NWWWNWWN. There are $\binom{8}{3}$ different ways to choose where the Ns go in the string, after which the positions of the Ws is completely determined. Alternatively, we can pick where the Ws go, in $\binom{8}{5}$ ways. Thus the number of possible paths from the café to the sandwich shop is

$$\binom{8}{3} = \frac{8!}{3!5!} = \binom{8}{5} = 56.$$

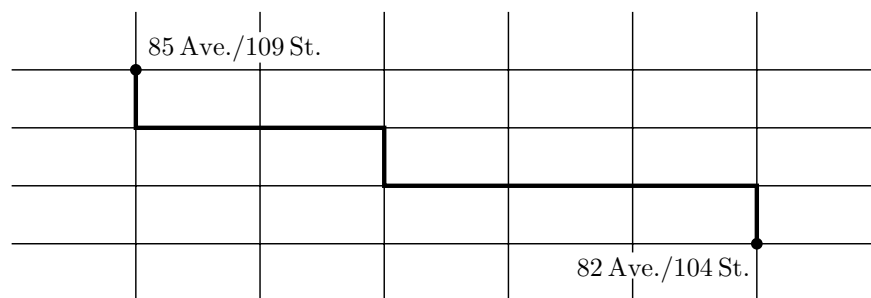


Fig. 12. One possible path from the café to the sandwich shop.

Here we see an important principle in action. *The number of k -combinations of a set of size n is the same as the number of $(n - k)$ -combinations of the same set.* This is because choosing k items is the same as choosing which $n - k$ items *not* to take. In general, the number of ways of walking m blocks in one direction and n blocks in a perpendicular direction is

$$\binom{m+n}{m} = \binom{m+n}{n}.$$

Stars and bars. Suppose there is a shop selling marbles that come in six colours. How many ways are there to select a bag of four marbles? Imagine that the marbles are laid out in six different bins all in a row, one for each colour. We stand at the first bin and may take some number of them. Mark down a star for each one that we take. Then we step forward to the next bin; record this action by drawing a bar. If we take some marbles here, mark down a star for each marble taken. Then we step forward again, so mark down a bar. At the end, we must have taken exactly four marbles and switched bins five times, so there must be exactly four stars and five bars, e.g., the sequence

$$| \quad | \star | \star | \quad | \star \star$$

corresponds to taking one marble of the third colour, one of the fourth, and two of the sixth. There are nine slots, four of which have to be stars, so the number of ways to pick four marbles out of six bins is $\binom{9}{4} = 126$. In general the number of ways to pick m marbles out of b bins is $\binom{m+b-1}{m}$. This method of counting is called the *stars-and-bars* technique.

A night at the casino. Suppose we are playing a high-stakes card game at the local gambling house and it is somehow pertinent for us to calculate how many five-card hands there are with at least two hearts. A naïve attempt would run as follows. We pick two cards to be hearts, in $\binom{13}{2}$ ways, then there are 50 cards left to choose from, and we can pick any three of those, in $\binom{50}{3}$ ways.

This seems logical, except that we have counted some hands twice! For instance, suppose that in the first step we pick our two hearts to be $2\heartsuit$ and $6\heartsuit$, then for our remaining three cards we pick $8\heartsuit$, $3\diamondsuit$, and $Q\spadesuit$. We can get the exact same hand by first picking our two hearts to be $6\heartsuit$ and $8\heartsuit$, and then choosing $2\heartsuit$, $3\diamondsuit$, and $Q\spadesuit$ for the remaining three cards.

One correct way of getting the correct count is to count the number of hands with exactly two hearts, exactly three hearts, exactly four hearts, and exactly five hearts. Since these possibilities all define disjoint sets, we can simply add them up at the end. The number of ways of drawing exactly two hearts is

$$\binom{13}{2}\binom{52-13}{3} = \binom{13}{2}\binom{39}{3} = 712842,$$

since there are 39 non-heart cards and we must select three of them. Similarly, the number of ways of drawing exactly three hearts is

$$\binom{13}{3}\binom{52-13}{2} = \binom{13}{3}\binom{39}{2} = 211926;$$

we also calculate

$$\binom{13}{4}\binom{39}{1} = 27885 \quad \text{and} \quad \binom{13}{5}\binom{39}{0} = 1287.$$

Hence the number of five-card hands with at least two hearts is

$$712842 + 211926 + 27885 + 1287 = 953940.$$

This number has little practical value unless we know the number of five-card hands in total; this is $\binom{52}{5} = 2598960$. We can then calculate the probability of a five-card hand containing at least two hearts to be $953940/2598960 \approx 36.7\%$.

Deriving the total number of five-card hands suggests a different way to get the same count. We could take this total and subtract the number of hands with no hearts, as well as the number of hands with exactly one heart. This gives

$$\binom{52}{5} - \binom{13}{0}\binom{39}{5} - \binom{13}{1}\binom{39}{4} = 2598960 - 575757 - 1069263 = 953940,$$

the same figure we obtained in the previous paragraph.

Identities involving binomial coefficients. We already observed earlier that the number of ways of picking k objects out of a set of size n is the same as the number of picking $n-k$ objects out of the same set. This can be seen algebraically by expanding

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} = \frac{n!}{(n-k)!(n-(n-k))!} = \binom{n}{n-k}.$$

Here is another identity involving binomial coefficients.

Proposition 87. *Let $n \geq 0$ be an integer. Then*

$$\sum_{k=0}^n \binom{n}{k}^2 = \binom{2n}{n}.$$

Proof. Let X be set of size $2n$. The right-hand side counts the number of subsets of X of size n . We count this in a different way. Let A and B be such that $|A| = |B| = n$, and $A \cup B = X$ (so we must have $A \cap B = \emptyset$). Then choosing a subset of X is the same as choosing k elements of A , where $0 \leq k \leq n$, and then choosing the remaining $n - k$ elements from B . In other words,

$$\binom{2n}{n} = \sum_{k=0}^n \binom{n}{k} \binom{n}{n-k} = \binom{n}{k}^2. \quad \blacksquare$$

More identities can be obtained by laying out binomial coefficients in the triangle

$$\begin{array}{ccccccc} & & & \binom{0}{0} & & & \\ & & \binom{1}{0} & & \binom{1}{1} & & \\ & \binom{2}{0} & & \binom{2}{1} & & \binom{2}{2} & \\ \binom{3}{0} & & \binom{3}{1} & & \binom{3}{2} & & \binom{3}{3} \\ \vdots & & \vdots & & \vdots & & \vdots \end{array}$$

Adding a few more rows and writing out the actual values of the binomial coefficients yields Fig. 13. This triangle is named after B. Pascal, who wrote a treatise regarding some of its properties that was published posthumously in 1665. (The triangle is known by different names in some different countries, as it was studied by Al-Karaji in the tenth or eleventh century, by O. Khayyám about a hundred years later, by Yang Hui in the thirteenth century, and by Tartaglia in 1556, among many others.)

$$\begin{array}{cccccccccccccccc} & & & & & & & 1 & & & & & & & & & \\ & & & & & & 1 & & 1 & & & & & & & & \\ & & & & & 1 & & 2 & & 1 & & & & & & & \\ & & & 1 & & 3 & & 3 & & 1 & & & & & & & \\ & & 1 & & 4 & & 6 & & 4 & & 1 & & & & & & \\ & & 1 & & 5 & & 10 & & 10 & & 5 & & 1 & & & & \\ & 1 & & 6 & & 15 & & 20 & & 15 & & 6 & & 1 & & & \\ & 1 & & 7 & & 21 & & 35 & & 35 & & 21 & & 7 & & 1 & \\ 1 & & 8 & & 28 & & 56 & & 70 & & 56 & & 28 & & 8 & & 1 \\ \vdots & & \vdots & & \vdots & & \vdots & & \vdots & & \vdots & & \vdots & & \vdots & & \vdots \end{array}$$

Fig. 13. The first nine rows of Pascal's triangle

This triangle can easily be drawn by hand by writing 1s along its boundary, and then filling it using the observation that every number is the sum of the two directly above it in the previous row. Of course, this rule has to be proved, so we shall do so right now.

Theorem 88 (*Pascal's identity*). *Let $n \geq 0$ and $k \geq 1$ be integers. Then*

$$\binom{n}{k} = \binom{n-1}{k} + \binom{n-1}{k-1}.$$

Proof. Set $X = \{1, 2, \dots, n\}$. We count the number of sets $Y \subseteq X$ there are of size k . Either $1 \notin Y$ or $1 \in Y$. If $1 \notin Y$, then there are $\binom{n-1}{k}$ possibilities for what Y can be. If $1 \in Y$, then there are still $k-1$ other elements of Y , and $n-1$ elements to choose from, so there are $\binom{n-1}{k-1}$ possibilities in this case. Hence the right-hand side counts the number of subsets of X with cardinality k . But this is exactly what the left-hand side counts as well. ■

The only number that appears infinitely many times is 1, since any $n \in \mathbf{Z}$ can only appear in the first $n+1$ rows of the triangle (if $m > n$, then for all $0 \leq k \leq m$, $\binom{m}{k}$ is either 1 or at least m). A fun conjecture to think about is Singmaster's conjecture, put forth by D. Singmaster in 1971: *There exists some number C such that every integer $n \geq 2$ appears in Pascal's triangle at most C times.* Whether this statement is true is still an open problem. The integer (besides 1) known to appear most often in Pascal's triangle is 3003. It appears eight times, and no other integer is known to appear eight times.

But let us not be too diverted by lofty conjectures. Returning to things we can prove, another thing we notice about Pascal's triangle is that the first row sums to 1, the second row sums to 2, the third to 4, the fourth to 8, and so on. We might therefore conjecture that the n th row of Pascal's triangle should always sum up to 2^n . Indeed, we have the following easy corollary of the binomial theorem.

Proposition 89. *For all $n \in \mathbf{N}$,*

$$\sum_{k=0}^n \binom{n}{k} = 2^n.$$

Proof. By the binomial theorem with $x = y = 1$, we have

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

The proposition can also be proven by noting that the left-hand side counts the number of subsets of $\{1, 2, \dots, n\}$ of size 0, plus the number of subsets of size 1, and so on up to all subsets of size n . But adding these all up simply gives the total number of subsets of $\{1, 2, \dots, n\}$, which we know to be 2^n . ■

The binomial theorem with a different choice of x and y allows us to reprove a statement we showed a long time ago, that for any finite nonempty set X , the number of subsets of X with odd cardinality is the same as the number of subsets of X with even cardinality.

Another proof of Proposition 23. We re-establish the notation of the proposition statement. Let X be a finite nonempty set, let E be the set of subsets of X with even cardinality, and let O be the set of subsets with odd cardinality. By the binomial theorem with $x = -1$ and $y = 1$, we have

$$0 = (-1 + 1)^n = \sum_{k=0}^n \binom{n}{k} (-1)^k 1^{n-k} = \binom{n}{0} - \binom{n}{1} + \binom{n}{2} - \cdots + (-1)^n \binom{n}{n}.$$

Hence

$$\binom{n}{0} + \binom{n}{2} + \cdots = \binom{n}{1} + \binom{n}{3} + \cdots,$$

but the left-hand side of this is just $|E|$, and the right-hand side is $|O|$. ■

The freshman's dream. Back when you were in high school, you might have made the following computational error:

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

It turns out you were not wrong, provided x and y were integers and you were working modulo 2. In general, one has the following consequence of the binomial theorem.

Theorem 90 (*Freshman's dream*). *Let p be a prime number and let $x, y \in \mathbf{Z}$. Then*

$$(x + y)^p \equiv x^p + y^p \pmod{p}.$$

Proof. By the binomial theorem,

$$(x + y)^p = \sum_{k=0}^p \binom{p}{k} x^k y^{p-k} = x^p + y^p + \sum_{k=1}^{p-1} \binom{p}{k} x^k y^{p-k}.$$

Taking this equation modulo p , we are done if we can show that for all $1 \leq k \leq p-1$, $\binom{p}{k} \equiv 0 \pmod{p}$.

Let $1 \leq k \leq p-1$ and expand

$$\binom{p}{k} = \frac{p!}{k!(p-k)!}$$

Multiplying both sides by $k!(p-k)!$, we get

$$k!(p-k)! \binom{p}{k} = p! = p(p-1)!.$$

This is an integer and the right-hand side is a multiple of p , so the left-hand side is as well. Then since p divides the product $k!(p-k)!\binom{p}{k}$, either p divides $k!(p-k)!$ or it divides $\binom{p}{k}$. But

$$k!(p-k)! = k(k-1)\cdots 2\cdot 1\cdot (p-k)(p-k-1)\cdots 2\cdot 1,$$

and all of the factors on the right-hand side are positive integers less than p . So p cannot divide any of them, and hence p does not divide $k!(p-k)!$. We conclude that $p \mid \binom{p}{k}$. ■

26.XI **A day spent selling milk.** The binomial theorem can be used to find the value of a specific coefficient in a polynomial, without necessarily expanding the whole thing out. To illustrate this, consider the following scenario. Imagine that a farmer runs a roadside milk stand. Each time a customer visits they either

- i) pay her \$6 and buy one bottle of white, chocolate, or strawberry milk; or
- ii) sell her back a glass bottle for \$1.

So after one customer visits, there are three ways in which the farmer could have made \$6, and one way in which she could have made $-\$1$. Let us arrange these possibilities as a function of x ; we write $3x^6 + x^{-1}$.

What benefit is there of doing this? Well, suppose two customers come by. There are three ways for the first customer to give her \$6, and three ways for the second customer to give her \$6, so the number of ways for her to make \$12 after two customers visit is 9 (where the order of the flavours she sells matters—selling chocolate then strawberry is different from selling strawberry then chocolate). It is also possible for her to net \$5 from two customers arriving. She could buy a bottle from the first customer and then sells the second a bottle of milk, or vice versa. There are three ways for each possibility, so six ways in total that she can make \$5. Lastly, there is only one way for her to have made $-\$2$. This is if she buys back a bottle from both of them. Thus the function of x that describes all the fiscal possibilities of the farmer is

$$9x^{12} + 6x^5 + x^{-2} = (3x^6 + x^{-1})^2.$$

By a similar reasoning, the number of ways for the farmer to make x dollars after n customers come by is $(3x^6 + x^{-1})^n$. When n gets large we don't necessarily want to expand out this entire function (it has 2^n terms), but the binomial theorem allows this representation to be useful nonetheless.

Suppose the farmer expects 15 customers to come by tomorrow. Let's calculate the number of ways that the farmer can make \$48 from this. By the

binomial theorem we have

$$\begin{aligned}
 (3x^6 + x^{-1})^{15} &= \sum_{k=0}^{15} \binom{15}{k} (3x^6)^k (x^{-1})^{15-k} \\
 &= \sum_{k=0}^{15} \binom{15}{k} 3^k x^{6k-15+k} \\
 &= \sum_{k=0}^{15} \binom{15}{k} 3^k x^{7k-15},
 \end{aligned}$$

and the number of ways to make \$48 from 15 customers is the coefficient of x^{48} on the right-hand side. So we want to know if $7k - 15$ ever equals 48 for any integers $0 \leq k \leq 15$; setting $k = 9$ does the trick. Hence on the right-hand side, when $k = 9$, we have the term $\binom{15}{9} 3^9 x^{48}$, meaning that the number of ways of making \$48 after 15 customers come by is $\binom{15}{9} 3^9 = 98513415$.

A sum of powers of x whose coefficients encode some combinatorial meaning (the number of ways to do something) is called a *generating function*. Using generating functions is one of the most beautiful and powerful techniques in combinatorics; the topic is usually covered in a second course on discrete mathematics (e.g., MATH 340).

*A generating function is a clothesline
on which we hang up a sequence for display.*

— HERBERT S. WILF, *generatingfunctionology* (1989)

19. Applications of the pigeonhole principle

Back in Section 5, we proved the pigeonhole principle: *If n pigeons nest in $n - 1$ holes, then at least one hole contains more than one pigeon.* We now present some examples illustrating its use.

Sheep. Here is the first example. Suppose five sheep are placed in a square enclosure with a side length of 20 metres. One can prove that any point in time, there are always two sheep within a distance of $\sqrt{200} \approx 14.14$ metres from one another.

To use the pigeonhole principle, we need to decide what will be our “pigeons” and what will be our “holes”. In this case, we can divide the enclosure into four square quadrants, each of side length 10. We take our “pigeons” to be the sheep and our “holes” to be the quadrants. Since there are five sheep and four quadrants, there is at least one quadrant that contains more than one sheep. But the quadrants are squares with side lengths of 10 metres, so the farthest any two sheep can stand apart from each other while being in the same quadrant is the length of the diagonal. This is $\sqrt{10^2 + 10^2} = \sqrt{200}$ by the Pythagorean theorem.

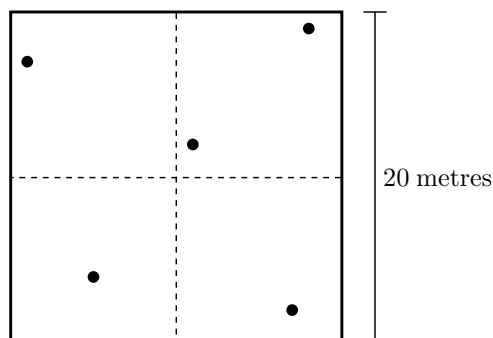


Fig. 14. Five sheep in an enclosure.

Subset sums. Let $A \subseteq \{1, 2, \dots, 100\}$ with $|A| = 10$. Prove that there has to be subsets $X \subseteq A$ and $Y \subseteq A$ with $X \neq Y$ such that

$$\sum_{x \in X} x = \sum_{y \in Y} y.$$

We define $f : 2^A \rightarrow \mathbf{N}$ by letting

$$f(X) = \sum_{x \in X} x.$$

(So $f(\emptyset) = 0$.) What is the maximum possible sum? Well, this happens when $A = \{100, 99, 98, \dots, 91\}$ and $X = A$, in which case

$$f(X) = \sum_{x \in X} x \leq 100 + 99 + 98 + \dots + 91 \leq 1000.$$

The minimum possible sum is 0, when $X = \emptyset$. Hence let 2^A be the set of pigeons and let $\{0, 1, 2, \dots, 1000\}$ be the set of holes. There are $2^{|A|} = 2^{10} = 1024$ pigeons and 1001 holes, so by the pigeonhole principle there must be two pigeons that nest in the same hole; in other words, f is not injective. Hence there are $X, Y \in 2^A$ with $X \neq Y$ such that $f(X) = f(Y)$.

Theoretical limits to data compression. For $n \in \mathbf{N}$, let B_n be the set of all binary strings of length n . Let B be the set of all binary strings of any finite length, so that

$$B = B_0 \cup B_1 \cup B_2 \cup \dots = \bigcup_{n=0}^{\infty} B_n.$$

Consider a string $x \in B$ as a file (containing some data encoded by the 0s and 1s). Let $L(x)$ be the length of x , i.e., the number of bits needed to encode x . Consider an algorithm α that takes a file as input and outputs a file. We can view α as a function $\alpha : B \rightarrow B$. Now imagine that the goal of α is to *compress*

the file x ; that is, we would ideally like to feed a file x to α and get a file y that is shorter in length. A compression algorithm α is *lossless* that for any y in the range of α , we can recover x such that $\alpha(x) = y$. In other words, α is *lossless* if and only if it is injective when considered as a function $\alpha : B \rightarrow B$. Examples of widely-used file formats that feature lossless compression algorithms are ZIP, PNG, and GIF.

What we will prove is that any compression algorithm that decreases the file size for some inputs must increase the file size for other inputs. (For example, there exist some files such that compressing them into a ZIP file causes them to get larger!) This is formalised in the following theorem.

Theorem 91. *Let $\alpha : B \rightarrow B$ be a lossless compression algorithm such that for at least one $x \in B$, we have $L(\alpha(x)) < L(x)$. Then there exists $y \in B$ such that $L(\alpha(y)) > L(y)$.*

Proof. For a contradiction, assume that there exists $\alpha : B \rightarrow B$ that is injective such that for at least some $x \in B$, $L(\alpha(x)) < L(x)$, and for all $y \in B$, $L(\alpha(y)) \leq L(y)$.

Pick $x \in B$ to be of minimal length such that $L(\alpha(x)) < L(x)$. Set $n = L(x)$ and $m = L(\alpha(x))$, so that $n > m$. By minimality of x , if $y \in B$ with $L(y) < n$, then $L(\alpha(y)) \geq L(y)$, which, combined with our earlier assumption, means that $L(\alpha(y)) = L(y)$ for all $y \in B$ with $L(y) < n$.

Define a new algorithm that has a restricted domain and codomain. We let $\alpha' : (B_m \cup \{x\}) \rightarrow B_m$ be given by $\alpha'(x) = \alpha$. The codomain is correct because $\alpha(y) \in B_m$ for all $y \in B_m$ and $\alpha(x) \in B_m$ as well, by choice of x . Let $B_m \cup \{x\}$ be the set of pigeons and let B_m be the set of pigeonholes. Since there are $|B_m| + 1$ pigeons and $|B_m|$ holes, there must be some $y, y' \in B_m \cup \{x\}$ with $y \neq y'$ such that $\alpha'(y) = \alpha'(y')$. But this means that $\alpha(y) = \alpha(y')$ as well. Hence α is not injective: a contradiction. ■

20. Recurrences

Let $(a_n) = a_0, a_1, a_2, \dots$ be a sequence of real numbers (in all the examples we consider, they will be integers). A *recurrence relation* for the sequence is a formula for a_n in terms of some of the a_m with $m < n$. For instance, the *Fibonacci sequence* is the sequence F_n defined by $F_0 = 0$, $F_1 = 1$, and $F_n = F_{n-1} + F_{n-2}$ for $n \geq 2$. The first few terms of F_n are

$$0, 1, 1, 2, 3, 5, 8, 13, 21, \dots$$

The sequence is named for the Italian mathematician Fibonacci, who in his 1202 book *Liber Abaci* used these numbers to describe the growth of a population of rabbits under certain idealised conditions. The sequence was known to Indian linguists many centuries earlier, as it pertains to the number of different ways to arrange long and short Sanskrit syllables into lines of poetry.

28.XI **Finding recurrence relations.** First we show an example of how to find a recurrence relation for a given sequence. For $n \geq 1$, let C_n be the number of binary strings of length n not containing 00. We have $C_1 = 2$ (0 and 1 are both fine) and $C_2 = 3$ (01, 10, and 11 are all fine, but 00 is not). Now let $n \geq 2$ and consider how to derive a count for C_n , assuming we know the smaller values of n . Consider a string of length n with no 00. Either the string ends in 1 or it ends in 10, since if it ends with a 0 the second-last bit must be 1. In the first case, there are $n - 1$ bits preceding the 1, and the only requirement is that this string of $n - 1$ bits not contain 00. So there are C_{n-1} possibilities in the first case. In the second case, there are $n - 2$ bits preceding the 10, and the requirement again is that the string $n - 2$ bits not contain 00. Hence in the second case there are C_{n-2} possibilities. Adding up the two possibilities, we have

$$C_n = C_{n-1} + C_{n-2}.$$

This is the exact same recurrence relation as for the Fibonacci numbers! The only difference is that the answer is shifted (e.g., $C_1 = 3 = F_3$). So in general, $C_n = F_{n+2}$ for all integers $n \geq 1$.

Why don't we do another one. Let A_n be the number of strings of length n using the symbols 0, 1, 2, 3, such that there are an even number of 1s. We have $A_1 = 3$, and we can calculate $A_2 = 10$, since there are 16 possible strings, and the forbidden ones are 12, 13, 14, 21, 31, and 41. Now let $n \geq 2$ and consider a string of length n with an even number of 1s. Either the last number is 1 or it is not. In the first case, there are $n - 1$ bits preceding the 1, and we need them to have an *odd* number of bits. The number of ways of doing this is $4^{n-1} - A_{n-1}$ (subtract the ones with an even number of 1s from the total). In the second case, the last digit is 0, 2, or 3, and in each case we just need the preceding $n - 1$ bits to have an even number of 1s, in A_n ways. Hence there are $3A_{n-1}$ possibilities for this case. Adding everything up, we have

$$A_n = 3A_{n-1} + 4^{n-1} - A_{n-1} = 2A_{n-1} + 4^{n-1}.$$

Solving recurrences. Suppose we have a recurrence relation for a sequence a_n . We would like to “solve” this recurrence by deriving a formula for a_n that isn't recursive, that is, one that does not contain any instances of a_m for $m < n$ on the right-hand side.

For example, in the very simple case where the recurrence has only one term, we have the following solution.

Proposition 92. Suppose (a_n) is a sequence with $a_0 = c$ for some $c \in \mathbf{R}$ and, for all $n \geq 0$, $a_n = b \cdot a_{n-1}$ for some $b \in \mathbf{R}$. Then for all $n \geq 0$, we have

$$a_n = cb^n.$$

Proof. By induction on n . For $n = 0$, we have

$$a_0 = c = cb^0.$$

(We define $0^0 = 1$). Now suppose the formula proven for n and observe that

$$a_{n+1} = b \cdot a_n = b \cdot cb^n = cb^{n+1}. \quad \blacksquare$$

That was too easy. In general, we will consider recurrences of the form

$$a_n = f_1(n)a_{n-1} + f_2(n)a_{n-2} + \cdots + f_k(n)a_{n-k} + g(n)$$

where $k \geq 1$ is an integer and $f_1, f_2, \dots, f_k, g : \mathbf{N} \rightarrow \mathbf{R}$ are functions. These recurrences are said to be *linear of degree k* . If $g(n) = 0$, the recurrence is called *homogeneous*; otherwise, it is called *non-homogeneous*.

For a non-homogeneous recurrence, the same recurrence without the $g(n)$ term is called the *associated homogeneous recurrence*. If, for all $1 \leq i \leq k$, the function $f_i(n)$ is a constant function, then the recurrence is said to have *constant coefficients*.

A sequence (p_n) that satisfies the recurrence is called a *particular solution*. A *general solution* is a formula describing all possible solutions using some parameters. If we specify values for the first few terms a_0, a_1, a_2, \dots , these are called *initial conditions* for the recurrence.

Non-homogeneous recurrences of degree 1. So far, the only recurrences we know how to solve are linear homogeneous recurrences of degree 1. As a step up in difficulty, let us now consider non-homogeneous recurrences of degree 1.

Theorem 93. *Consider the non-homogeneous recurrence*

$$a_n = f(n)a_{n-1} + g(n),$$

where $f, g : \mathbf{N} \rightarrow \mathbf{R}$. If (p_n) is any particular solution to the recurrence and (h_n) is a general solution to the associated homogenous recurrence, i.e.,

$$h_n = f(n)h_{n-1}$$

for all $n \geq 1$, then the general solution for the recurrence is

$$a_n = h_n + p_n.$$

Proof. Let (b_n) be any solution to the recurrence and let (p_n) be a particular solution to the recurrence, so that for $n \geq 1$ we have

$$b_n = f(n)b_{n-1} + g(n) \quad \text{and} \quad p_n = f(n)p_{n-1} + g(n).$$

Let (c_n) be the sequence given by $c_n = b_n - p_n$. Note that for all $n \geq 1$,

$$\begin{aligned} f(n)c_{n-1} &= f(n)(b_{n-1} - p_{n-1}) \\ &= f(n)b_{n-1} - f(n)p_{n-1} \\ &= b_n - g(n) - (p_n - g(n)) \\ &= b_n - p_n \\ &= c_n. \end{aligned}$$

So (c_n) is a solution to the associated homogeneous recurrence. In other words, it is obtained from (h_n) by specifying values for its parameters. Hence, letting

$$a_n = h_n + p_n,$$

the solution $b_n = c_n + p_n$ is derived by specifying specific values for parameters in the general solution (a_n) . ■

Here, then, is the general method for solving a non-homogeneous recurrence of degree 1:

- i) Solve the associated homogeneous recurrence, using Proposition 92, to get a general solution (h_n) with parameters.
- ii) Find a solution p_n to the recurrence. To do so, we often find it useful to guess that p_n should have a similar form to $g(n)$. We'll see an example of this soon.
- iii) Theorem 93 tells us that the general solution to our recurrence is

$$a_n = h_n + p_n.$$

- iv) There will be unspecified parameters in the general solution, so if possible, we solve for these parameters by calculating some initial conditions.

To illustrate this method, let us revisit the example from earlier about the number A_n of strings of 0, 1, 2, and 3 that have an even number of 1s. We derived the recurrence

$$A_n = 2A_{n-1} + 4^{n-1}.$$

Now let's solve it using Proposition 92 and Theorem 93. First, we find the general solution to the associated homogeneous recurrence $A_n = 2A_{n-1}$. By Proposition 92, this is

$$h_n = \alpha 2^n,$$

where $\alpha \in \mathbf{R}$ is a parameter that we can fill in by looking at initial conditions, but we'll skip that for now. Now we want a particular solution p_n . Since $g(n) = 4^{n-1}$, we guess that $p_n = \beta 4^n$ for some $\beta \in \mathbf{R}$. We want $p_n = 2p_{n-1} + 4^{n-1}$, so

$$\beta 4^n = 2\beta 4^{n-1} + 4^{n-1}.$$

Dividing through by 4^{n-1} yields

$$4\beta = 2\beta + 1,$$

so we can take $\beta = 1/2$, and $p_n = (1/2)4^n$ is a particular solution. By Theorem 93, the general solution to our original recurrence is

$$\alpha \cdot 2^n + \frac{1}{2}4^n.$$

If $g(n)$ is then take p_n to be ...
a constant q ,	a constant r .
a linear function $q_0 + q_1n$,	a linear function $r_0 + r_1n$.
an exponential function qt^n ,	an exponential function with the same base rt^n .

Table 1. Particular solutions p_n for certain classes of $g(n)$.

Now, earlier we derived the initial condition $A_1 = 3$, so

$$3 = A_1 = \alpha \cdot 2^1 + \frac{1}{2}4^1.$$

Solving this, we obtain $\alpha = \frac{1}{2}$. So for $n \geq 1$, the final formula for A_n is

$$A_n = \frac{1}{2}2^n + \frac{1}{2}4^n.$$

We verify that

$$A_2 = \frac{1}{2}2^2 + \frac{1}{2}4^2 = 2 + 8 = 10,$$

which matches the count we performed earlier.

In this example, we guessed that the particular solution p_n should take the form $\beta 4^n$ for some β , since $g(n) = 4^{n-1}$. This sort of guess works for certain other classes of $g(n)$ as well. The general principle is outlined in Table 1.

Theorem 93 is a special case of a more general theorem, whose proof is beyond the scope of this course.

Theorem 94. *Consider the recurrence*

$$a_n = f_1(n)a_{n-1} + f_2(n)a_{n-2} + \cdots + f_k(n)a_{n-k} + g(n)$$

for some $k \geq 1$ and $f_1, f_2, \dots, f_k, g : \mathbf{N} \rightarrow \mathbf{R}$. If (p_n) is any particular solution and (h_n) is a solution to the associated homogeneous recurrence, i.e.,

$$h_n = f_1(n)a_{n-1} + f_2(n)a_{n-2} + \cdots + f_k(n)a_{n-k}$$

for all $n \geq 1$, then the general solution for (a_n) is given by $a_n = h_n + p_n$. **■**

03.XII **Recurrences with constant coefficients.** We now turn to the case where the f_i are all constant functions. This means there are constants c_1, c_2, \dots, c_k such that $f_1(n) = c_1$ for all n , $f_2(n) = c_2$ for all n , and so on. In light of Theorem 94, we can restrict our attention to homogeneous recurrences, since we can solve

non-homogeneous recurrences by solving their associated homogeneous recurrences and then applying the theorem. Homogeneous recurrences with constant coefficients look like

$$a_n = c_1 a_{n-1} + c_2 a_{n-2} + \cdots + c_k a_{n-k},$$

for some degree $k \geq 1$. We define the *characteristic equation* of such a recurrence to be the equation

$$x^k = c_1 x^{k-1} + c_2 x^{k-2} + \cdots + c_k,$$

and the *characteristic polynomial* is the polynomial obtained by subtracting the right-hand side from the left-hand side:

$$x^k - c_1 x^{k-1} + c_2 x^{k-2} + \cdots + c_k$$

It turns out that this polynomial entirely determines the general solution to the recurrence. We shall state and prove this theorem in the special case that $k = 2$ and the two roots of the characteristic polynomial are distinct.

Theorem 95. *Consider the recurrence given by*

$$a_n = c_1 a_{n-1} + c_2 a_{n-2}.$$

If the characteristic polynomial of this recurrence has two distinct roots r_1 and r_2 , then the general solution of the recurrence is

$$a_n = \alpha_1 r_1^n + \alpha_2 r_2^n.$$

Supposing we know the initial conditions a_0 and a_1 , we have the identities

$$\alpha_1 = \frac{a_0 r_2 - a_1}{r_2 - r_1} \quad \text{and} \quad \alpha_2 = \frac{a_1 - a_0 r_1}{r_2 - r_1}.$$

Proof. The characteristic polynomial of the recurrence is

$$x^2 - c_1 x + c_2.$$

For any root r of this polynomial, we have

$$r^2 = c_1 r + c_2,$$

so for any $n \geq 2$ we can multiply both sides by r^{n-2} to get

$$r^n = c_1 r^{n-1} + c_2 r^{n-2}.$$

Hence the sequence (r^n) is a solution to the recurrence. Conversely, if r is such that (r^n) is a solution to the recurrence, then $r^n = c_1 r^{n-1} + c_2 r^{n-2}$ for all

$n \geq 2$, so we can divide by r^{n-2} to get $r^2 = c_1r + c_2$; that is, r is a root of the characteristic polynomial.

What we have found so far is that $a_n = r_1^n$ and $a_n = r_2^n$ are both solutions to the recurrence, and that they are the only two solutions that take the form r^n for some r . Now we must show that for any scalars α_1 and α_2 ,

$$f_n = \alpha_1 r_1^n + \alpha_2 r_2^n$$

is also a solution to the recurrence. Well,

$$\begin{aligned} c_1 f_{n-1} + c_2 f_{n-2} &= c_1 (\alpha_1 r_1^{n-1} + \alpha_2 r_2^{n-1}) + c_2 (\alpha_1 r_1^{n-2} + \alpha_2 r_2^{n-2}) \\ &= c_1 \alpha_1 r_1^{n-1} + c_1 \alpha_2 r_2^{n-1} + c_2 \alpha_1 r_1^{n-2} + c_2 \alpha_2 r_2^{n-2} \\ &= \alpha_1 (c_1 r_1^{n-1} + c_2 r_1^{n-2}) + \alpha_2 (c_1 r_2^{n-1} + c_2 r_2^{n-2}) \\ &= \alpha_1 r_1^n + \alpha_2 r_2^n \\ &= f_n, \end{aligned}$$

So the claimed expression is in fact the general solution to the recurrence.

Lastly, we must show that we can solve for α_1 and α_2 once we know the initial conditions a_0 and a_1 . We have

$$a_0 = \alpha_1 r_1^0 + \alpha_2 r_2^0 = \alpha_1 + \alpha_2$$

and

$$a_1 = \alpha_1 r_1 + \alpha_2 r_2.$$

We can arrange these equations into the matrix equation

$$\begin{pmatrix} 1 & 1 \\ r_1 & r_2 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} = \begin{pmatrix} a_0 \\ a_1 \end{pmatrix}.$$

This equation is solvable if and only if the determinant $r_2 - r_1$ of the matrix $\begin{pmatrix} 1 & 1 \\ r_1 & r_2 \end{pmatrix}$ is nonzero. But we assumed $r_1 \neq r_2$, so this is the case. Using the formula for the inverse of a 2×2 matrix, we find that the inverse of $\begin{pmatrix} 1 & 1 \\ r_1 & r_2 \end{pmatrix}$ is

$$\frac{1}{r_1 - r_2} \begin{pmatrix} r_2 & -1 \\ -r_1 & 1 \end{pmatrix},$$

whence

$$\begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} = \frac{1}{r_2 - r_1} \begin{pmatrix} r_2 & -1 \\ -r_1 & 1 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix}.$$

We conclude that

$$\alpha_1 = \frac{a_0 r_2 - a_1}{r_2 - r_1} \quad \text{and} \quad \alpha_2 = \frac{a_1 - a_0 r_1}{r_2 - r_1},$$

which is what we wanted to show. \blacksquare

Let us use this theorem to derive a non-recursive formula for the Fibonacci numbers. Recall that they satisfy the recurrence $F_n = F_{n-1} + F_{n-2}$, so the characteristic polynomial is

$$x^2 - x - 1.$$

By the quadratic formula, the two roots of this polynomial are

$$r_1 = \frac{1 - \sqrt{5}}{2} \quad \text{and} \quad r_2 = \frac{1 + \sqrt{5}}{2}.$$

The roots are distinct; in fact, we have

$$r_2 - r_1 = \frac{1 + \sqrt{5} - (1 - \sqrt{5})}{2} = \frac{2\sqrt{5}}{2} = \sqrt{5}.$$

In any case, Theorem 95 applies, and we have the general solution

$$F_n = \alpha_1 r_1^n + \alpha_2 r_2^n.$$

We have $F_0 = 0$ and $F_1 = 1$, so using the formulas given by the theorem, we calculate

$$\alpha_1 = \frac{0 \cdot r_2 - 1}{\sqrt{5}} = \frac{-1}{\sqrt{5}}$$

and

$$\alpha_2 = \frac{1 - 0 \cdot r_2}{\sqrt{5}} = \frac{1}{\sqrt{5}}.$$

So

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

From this formula, it is not even clear that $F_n \in \mathbf{N}$ for all n , though it must be since the first two terms are integer and each term is the sum of the two previous entries.

The quantity $r_1 = (1 + \sqrt{5})/2 \approx 1.618$ is called the *golden ratio*, and is often denoted φ . We have $r_2 = (1 - \sqrt{5})/2 \approx -0.618$, so $(-1/\sqrt{5})r_2^n$ is very small in absolute value for large n . This means that

$$F_n \sim \frac{1}{\sqrt{5}} \left(\frac{1 + \sqrt{5}}{2} \right)^n \approx \frac{1.61^n}{\sqrt{5}}.$$

In other words,

$$\lim_{n \rightarrow \infty} \frac{-1}{\sqrt{5}} \left(\frac{1 - \sqrt{5}}{2} \right)^n = 0,$$

so

$$\lim_{n \rightarrow \infty} \frac{F_{n+1}}{F_n} = \frac{1 + \sqrt{5}}{2}.$$

Incidentally, it is not just the first few decimal places in the approximations of r_1 and r_2 that match. As an exercise, prove that $r_2 = 1/r_1 = 1 - r_1$.

Theorem 95 can be generalised to higher degrees k . Its proof is morally the same as the one we gave for the case $k = 2$.

Theorem 96. *Let $k \geq 2$ be an integer and consider the recurrence given by*

$$a_n = c_1 a_{n-1} + c_2 a_{n-2} + \cdots + c_k a_{n-k}.$$

If the characteristic polynomial of this recurrence has k roots r_1, r_2, \dots, r_k that are pairwise distinct, then the general solution of the recurrence is

$$a_n = \alpha_1 r_1^n + \alpha_2 r_2^n + \cdots + \alpha_k r_k^n. \quad \blacksquare$$

Returning to the setting $k = 2$, the case where the characteristic polynomial has a repeated root is not covered by Theorem 95. Instead, we have the following theorem.

Theorem 97. *Consider the homogeneous recurrence*

$$a_n = c_1 a_{n-1} + c_2 a_{n-2}.$$

If the characteristic polynomial of this recurrence has a repeated root r , then the general solution of the recurrence is

$$a_n = \alpha_1 r^n + \alpha_2 n r^n.$$

Supposing we know the initial conditions a_0 and a_1 , we have the identities

$$\alpha_1 = a_0 \quad \text{and} \quad \alpha_2 = \frac{a_1 - a_0 r}{r}.$$

Proof. By the same logic as in the proof of Theorem 95, we know that r^n is a solution to the recurrence. We show that $b_n = n r^n$ is also a solution. Since r is a repeated root to the characteristic polynomial $x^2 - c_1 x - c_2$, we have

$$x^2 - c_1 x - c_2 = (x - r)^2 = x^2 - 2rx + r^2,$$

so by comparing coefficients we conclude that $c_1 = 2r$ and $c_2 = -r^2$. Now by expanding

$$\begin{aligned} c_1 b_{n-1} + c_2 b_{n-2} &= c_2(n-1)r^{n-1} + c_2(n-2)r^{n-2} \\ &= r^{n-2}(c_1(n-1)r + c_2(n-2)) \\ &= r^{n-2}(c_1 nr - c_1 r + c_2 n - 2c_2) \\ &= r^{n-2}(2r^2 n - 2r^2 - r^2 n + 2r^2) \\ &= nr^n \\ &= b_n \end{aligned}$$

we see that $b_n = n r^n$ satisfies the recurrence as well. By similar reasoning as in the proof of Theorem 95, any linear combination of the two solutions (r^n) and $(n r^n)$ is also a solution.

Now we prove the further claim. Suppose we know a_0 and a_1 . Then

$$a_0 = \alpha_1 r^0 + \alpha_2 0r^0 = \alpha_1,$$

and

$$a_1 = \alpha_1 r^1 + \alpha_2 r^1 = a_0 r + \alpha_2 r,$$

and hence

$$\alpha_1 = a_0 \quad \text{and} \quad \alpha_2 = \frac{a_1 - a_0 r}{r}. \quad \blacksquare$$

Let's illustrate this theorem with a final example. We solve the recurrence

$$a_n = 4a_{n-1} - 4a_{n-2} + n,$$

with the initial conditions $a_0 = 1$ and $a_1 = 3$. First we solve the homogeneous recurrence

$$h_n = 4h_{n-1} - 4h_{n-2}.$$

The characteristic polynomial is $x^2 - 4x + 4 = (x - 2)^2$, so by the above theorem,

$$h_n = \alpha_1 2^n + \alpha_2 n 2^n.$$

Now we need a particular solution p_n . We have $g(n) = n$, which is a linear function in n , so we guess that $p_n = cn + d$ for some real numbers c and d . Then from the identity $p_n = 4p_{n-1} - 4p_{n-2} + n$, we have

$$\begin{aligned} cn + d &= 4(c(n-1) + d) - 4(c(n-2) + d) + n \\ 0 &= -cn - d + 4cn - 4c + 4d - 4cn + 8c - 4d + n \\ 0 &= (1-c)n + (4c-d). \end{aligned}$$

So $1 - c = 0$ and $4c - d = 0$, whence $c = 1$ and $d = 4$. So $p_n = n + 4$. This gives us the general solution

$$a_n = \alpha_1 2^n + \alpha_2 n 2^n + n + 4.$$

We cannot use Theorem 97 to get α_1 and α_2 , since that theorem concerned homogeneous recurrences. Instead we directly compute

$$1 = a_0 = \alpha_1 2^0 + \beta \cdot 0 + 0 + 4 = \alpha_1 + 4,$$

and derive $\alpha_1 = -3$, then find that

$$3 = a_1 = \alpha_1 2 + \alpha_2 \cdot 2 + 1 + 4 = -6 + 2\alpha_2 + 4.$$

Hence $4 = 2\alpha_2$ and $\alpha_2 = 2$, which gives us a final solution of

$$a_n = 3 \cdot 2^n + 2n \cdot 2^n + n + 4$$

to our recurrence.