# MATH1001 Introduction to Number Theory

Dr NJ Wright

17 October, 2024

# Contents

# Introduction

Although this course is ostensibly about Number Theory, one of it's main aims is to introduce you to the concepts and ideas surrounding mathematical proofs. Virtually all of the important ideas in chapters 3 onwards, will be proved in detail, with nothing left to chance.

We begin in Section 1 with an introduction to Mathematical Logic and Proof, with the aim of introducing the language and basic concepts necessary for much of university level mathematics. We will try to be precise about what we can (and cannot) do in a proof, and crucially we will consider strategies for constructing a proof. The material is largely self-contained but see the books by Martin Liebek (Liebek, 2015), Daniel Velleman (Velleman, 2006) and Daniel Solow (Solow, 2013) for some further reading material.

The Number Theory material in Sections 3 - 7 are based heavily on the course textbook Elementary Number Theory by Jones and Jones (Jones and Jones, 2006). Our ultimate aim is to justify Euler's Theorem, Theorem 7.3, which will prove the main mathematical device needed to describe and justify the RSA encryption scheme we shall meet in Section 8.

We end the course in Chapter 8 with a brief look at some modern cryptography based on elementary Number Theory. In particular we aim to describe the workings of the RSA cryptographic system and the Diffie-Helman-Merkle key exchange system. (If we have time we shall cover a small amount of "extra" cryptography in Chapter 9 - however this will not be examinable!)

The book by Jones and Jones (Jones and Jones, 2006) is the principal source but a good alternative is that by Kenneth Rosen (Rosen, 2010).

With the online version of these notes, we have provided a number of Appendices, which we encourage the interested reader to consult. '

*Acknowledgements: Many thanks and much credit for these notes goes to Prof. Graham Niblo and Dr. Jim Renshaw. A large part of this course is based*

# Chapter 1

# Introduction to Proofs and Mathematical Logic

## 1.1 Introduction

In this course we will study numbers, more precisely we will study **integers**.

The integer number system has the property that we can *add* and *subtract* numbers, and can also *multiply* them, but we **cannot** in general *divide* integers. With this in mind we will try to avoid writing divisions in number theory: if $a$ and $b$ are integers for which $a$ is divisible by $b$, and we want to consider the ratio $x$ of the two numbers then instead of writing $x = a/b$, we would write $a = bx$. This avoids the possibility of writing down an expression like $2/3$ which is **not defined** in the context of the integers.

Our aim in the course is to **prove** results about the integers. In school you may have seen a few proofs but throughout university, especially in the area of pure mathematics, you will find a much greater emphasis on proof. In this module we will give an introduction to logic and proof techniques.

Why is it important to prove results? Here are some examples of results that we can prove in number theory.

- $p$ is a prime number if and only if $(p-1)! + 1$ is divisible by $p$.
- For any integers $a, b, c$ the equation

$$ax + by = c$$

5

has integer solutions if and only if $c$ is divisible by all common factors of $a$ and $b$.

- There are infinitely many integer solutions to the equation

$$a^2 + b^2 = c^2$$

where $a, b, c$ have no common factors.

- For $p$ a prime other than $2$, the equation

$$n^2 + 1 = kp$$

has integer solutions if and only if $p - 1$ is divisible by $4$.

These are examples of Theorems from number theory. A **theorem** is a statement which we know to be true because it has been proved.

Results such as these are true and may be useful, but are not obvious. Proving them allows us to be sure that the result is always true and, once we know this, we can then use these results as facts to help answer other questions.

On the other hand we might spot a pattern and think that a result is always true, but without a proof we cannot be sure. As an example, consider the numbers $2^p - 1$ where $p$ is a prime.

| $p$ | $2^p - 1$ |
|-----|-----------|
| 2   | 3         |
| 3   | 7         |
| 5   | 31        |
| 7   | 127       |
| 11  | 2047      |

The observation that 3, 7, 31 and 127 are themselves all prime numbers might lead us to think that if $p$ is a prime then $2^p$ must also be prime. It might even convince us that 2047 is also prime, after all there are no obvious factors: it is easy to see it is not divisible by 2, 3, 5, 7, 11, and computing a bit further it is also not divisible by 13, 17, 19. But 2047 is not prime (you don't have to go much further to find a factor).

Consider this as a warning: Looking for patterns can be very valuable and instructive, but without a proof you cannot be sure if the pattern continues.

## 1.2 Statements, predicates and connectives

What is a statement?

A *logical statement* or *proposition* is an expression that can meaningfully be assigned the values of *true* or *false*. A statement might be written in words or mathematical symbols, or often as a combination of both.

A statement must be **true** or **false** but **not both**. It cannot be *half-true* or *sort-of-true*.

- "$1 > 0$" is a true statement
- "2047 is a prime number" is a false statement
- "What is the meaning of life?" is not a logical statement
- "if $n^2 = 4$ then what is $n$?" is a precise mathematical question (though with more than one answer) but is not a statement
- "$n^2 = 4$ has two integer solutions" is a true statement

A *predicate* is a 'statement' whose truth is dependent on one or more variables.

- "$n^2 = 4$" is an example of a *predicate*

Let us pause to consider what we mean by a *variable*. When you begin studying algebra you might think of variables as unknown numbers. But of course you often want to solve for the value of the variable at which point it is no longer unknown.

A variable also does not have to be a number. It might be a vector or a matrix: *let **v** denote the positive of the particle*. It could be a set: *let $S$ be the set of integer solutions of the equation $n^2 = 4$*. It could even by a function: *let $s$ denote the function whose value at a positive real number $y$ is the solution of the equation $x^2 = y$*.

Ultimately a *variable* is simply a placeholder which gives a name for some entity that you want to consider.

To consider the "algebra" of statements we will denote our statements by $P, Q, \ldots$

Note that we now have variables which are statements. Since statements must, by definition, have a well-defined (though not necessarily *known*) truth value, we might alternatively think of the variables $P, Q, \ldots$ as having values of true or false.

We can denote *predicates* by $P(n), P(x)$ etc. to indicate the variable on which they depend, e.g.

$$P(n): \quad n^2 = 4$$

## Logical Connectives

Many of the statements we meet in practice, are rather long and complicated statements. It is often useful to break them down into smaller parts and consider how these smaller parts contribute to the truth value of the original statement. We begin with the following three fundamental *logical connectives*:

1. *Conjunction*. The conjunction of two logical statements is true exactly when *both statements are true*. If $P$ and $Q$ are logical statements then the conjunction is denoted by $P \wedge Q$ (pronounced "$P$ and $Q$") and the truth value can be described using a *logic table* (or *truth table*).

$$
\begin{array}{cc|c}
P & Q & P \wedge Q \\
\hline
T & T & T \\
T & F & F \\
F & T & F \\
F & F & F
\end{array}
$$

2. *Disjunction*. The disjunction of two logical statements is true exactly when *at least one of the statements are true*. If $P$ and $Q$ are logical statements then the disjunction is denoted by $P \vee Q$ (pronounced "$P$ or $Q$") and the truth table is

$$
\begin{array}{cc|c}
P & Q & P \vee Q \\
\hline
T & T & T \\
T & F & T \\
F & T & T \\
F & F & F
\end{array}
$$

When we say "or" we will always mean the *inclusive or*, that is one or both of the statements are true unless we explicitly say that one or other but not both statements holds. For example when we say

*if $a$ is even or $b$ is even then $ab$ is even*

then the predicate "*$a$ is even or $b$ is even*" allows the cases - $a$ even and $b$ odd; - $a$ odd and $b$ even; and **also allows the case** - $a$ even and $b$ even.

3. *Negation*. The negation of a logical statement is true exactly when *the statement is false*. If $P$ is a logical statement then the negation is denoted by $\neg P$ (pronounced "not $P$") and the truth table is

$$
\begin{array}{c|c}
P & \neg P \\
\hline
T & F \\
F & T
\end{array}
$$

For example if $P$ is the statement "2047 *is a prime number*" then $\neg P$ means "*it is false that* 2047 *is a prime number*"; a statement which would normally be written / read as "2047 *is **not** a prime number*".

Note that bracketing of expressions is important: $(P \wedge Q) \vee R$ is different to $P \wedge (Q \vee R)$. Likewise $(\neg P) \wedge Q$ is different to $\neg(P \wedge Q)$.

By convention if we write $\neg P \wedge Q$ then we mean the former: $(\neg P) \wedge Q$. We will not write expressions such as $P \wedge Q \vee R$ as these are ambiguous.

When statements are written in words the "bracketing" may be ambiguous, so one must be careful to avoid confusion. How should the following statement be read?

*if $a$ is even or $b$ is even and $c$ is odd then $ab + c$ is odd*

From the context we might reasonably guess the intended meaning as being:

*if $((a$ is even $\vee \, b$ is even$) \wedge c$ is odd$)$ then $ab + c$ is odd*

To avoid ambiguity with writing this in words we might say something along the following lines:

*Suppose that $a$ is even and $b$ is even. Then if $c$ is odd then $ab + c$ is odd*

Conjunction, disjunction and negation are sufficient to build any other true/false-valued combinations of statements.

**Example 1.1.** Exclusive OR

| $P$ | $Q$ | $P$ xor $Q$ |
|:---:|:---:|:---:|
| $T$ | $T$ | $F$ |
| $T$ | $F$ | $T$ |
| $F$ | $T$ | $T$ |
| $F$ | $F$ | $F$ |

We can write $P$ xor $Q$ as

$$(P \vee Q) \wedge (\neg P \vee \neg Q).$$

Here $(P \vee Q)$ tells us that at least one of $P, Q$ is true while $(\neg P \vee \neg Q)$ says that at least one of $P, Q$ is false. Hence the conjunction of these two expressions will be true when exactly one of $P, Q$ is true and exactly one is false.

Another way to write $P$ xor $Q$ is as

$$(P \vee Q) \wedge \neg(P \wedge Q).$$

Again $(P \vee Q)$ tells us that at least one of $P, Q$ is true while $\neg(P \wedge Q)$ means "it is *not* true that both $P, Q$ are true" so again the conjunction tells us that exactly one of $P, Q$ is true.

Yet another way to write it is to say that we have two allowed cases: $P$ is true and $Q$ is false; $P$ is false and $Q$ is true. In symbols we can write this as the disjunction

$$(P \wedge \neg Q) \vee (\neg P \wedge Q).$$

Since each of these expressions are true in exactly the same cases we say that they are *logically equivalent*.

## 1.3 Logical equivalence

We say that two statements are *logically equivalent*, and denote this by $P \equiv Q$, if they have identical truth values. In other words $P$ is true exactly when $Q$ is true. So from our last example,

$$(P \vee Q) \wedge (\neg P \vee \neg Q) \equiv (P \vee Q) \wedge \neg(P \wedge Q) \equiv (P \wedge \neg Q) \vee (\neg P \wedge Q).$$

An expression which is always **true** independent of the truth of its components is called a *tautology*. For example

$$P \vee \neg P$$

is a tautology.

An expression which is always **false** independent of the truth of its components is called a *contradiction*. For example

$$P \wedge \neg P$$

is a contradiction.

The following theorem provides the algebraic rules for manipulating compound statements.

**Theorem 1.2**

Let $P, Q$ and $R$ be statements. Let $T$ and $F$ denote true and false.

1. Complement rules.
$$P \vee \neg P \equiv T$$
$$P \wedge \neg P \equiv F$$

2. Identity rules.
$$P \wedge T \equiv T \wedge P \equiv P$$
$$P \vee F \equiv F \vee P \equiv P$$

3. Idempotent rules.
$$P \wedge P \equiv P$$
$$P \vee P \equiv P$$

4. Commutative rules.
$$P \wedge Q \equiv Q \wedge P$$
$$P \vee Q \equiv Q \vee P$$

5. Associative rules.
$$(P \wedge Q) \wedge R \equiv P \wedge (Q \wedge R)$$
$$(P \vee Q) \vee R \equiv P \vee (Q \vee R)$$

6. Distributive rules.
$$(P \vee Q) \wedge R \equiv (P \wedge R) \vee (Q \wedge R)$$
$$(P \wedge Q) \vee R \equiv (P \vee R) \wedge (Q \vee R)$$

7. Double negation rule.
$$P \equiv \neg\neg P$$

8. De Morgan's rules.
$$\neg(P \wedge Q) \equiv (\neg P) \vee (\neg Q)$$
$$\neg(P \vee Q) \equiv (\neg P) \wedge (\neg Q)$$

The associative rules allow us to be a bit lazier with bracketing: we can write

$$P \wedge Q \wedge R$$

without brackets since it doesn't matter whether we interpret this as mean-

ing $(P \wedge Q) \wedge R$ or $P \wedge (Q \wedge R)$, and likewise for disjunction.

**Exercise 1.3.** As an exercise, construct the truth tables to justify each of these rules.

**Example 1.4.** As an example of how these rules can be used to deduce other rules we consider the *Absorption rules*:

$$P \vee (P \wedge Q) \equiv P$$

$$P \wedge (P \vee Q) \equiv P$$

Thinking about the first of these, the disjunction says that $P$ is true or $P \wedge Q$ is true (or as always both). If $P$ is true then, well, $P$ is true, while if $P \wedge Q$ is true then again $P$ must be true. So if the disjunction holds then $P$ must be true. But of course if $P$ is true then the disjunction is also true, which explains the equivalence.

We can alternatively show this using the logical rules from the theorem. As a first step let us deduce another equivalence (the relevance of which will become apparent in a moment). We will show the rule:

$$Q \vee T \equiv T \vee Q \equiv T$$

This is obvious from the truth-tables, but let us see how to deduce it via the logical equivalence rules. What do the rules tell us about $T$? The complement rule allows us to write

$$T \equiv Q \vee \neg Q$$

and therefore

$$Q \vee T \equiv Q \vee (Q \vee \neg Q) \equiv (Q \vee Q) \vee \neg Q$$

by the associative rule. The idempotent rule allows us to replace $Q \vee Q$ with $Q$ so we have

$$Q \vee T \equiv Q \vee \neg Q \equiv T$$

by again applying the complement rule. The same holds for $T \vee Q$ by commutativity.

With this new rule in our toolkit let us show the first absorption rule:

$$
\begin{aligned}
P \vee (P \wedge Q) &\equiv (P \wedge T) \vee (P \wedge Q) && \text{by the identity rule} \\
&\equiv P \wedge (T \vee Q) && \text{by the distributive rule} \\
&\equiv P \wedge T && \text{as } T \vee Q \equiv T \\
&\equiv P && \text{by the identity rule}
\end{aligned}
$$

We leave the second absorption rule as an exercise. As a hint, first consider the $F \wedge Q$.

**Example 1.5.** Here's an example of a word problem.

*There are three boxes on a table. One contains gold, the other two are empty. Each box has imprinted on it a clue to its contents. The clues are:*

> *Box A: "The gold is not here";*
>
> *Box B: "The gold is not here";*
>
> *Box C: "The gold is in Box B".*

*Only one of these clues is a true statement and the other two are false. By writing this information in symbolic form, determine which box has the gold?*

Let us write $A, B, C$ for the statements *"The gold is in box A"* etc. The clues then translate as:

> *1st clue:* $\neg A$
>
> *2nd clue:* $\neg B$
>
> *3rd clue:* $B$

How do we express the statement that *one of the clues is true and the other two are false*? This is a form of 3-way exclusive or: *1st is true and 2nd is false and 3rd is false, or the 1st is false and 2nd true and 3rd false, or 1st and 2nd false and 3rd true.*

Putting our clues into this the statement that exactly one clue is correct becomes:

$$(\neg A \wedge \neg\neg B \wedge \neg B) \vee (\neg\neg A \wedge \neg B \wedge \neg B) \vee (\neg\neg A \wedge \neg\neg B \wedge B)$$

Similarly the statement that there is gold in exactly one box is:

$$(A \wedge \neg B \wedge \neg C) \vee (\neg A \wedge B \wedge \neg C) \vee (\neg A \wedge \neg B \wedge C)$$

Taking the first statement we can begin simplifying using the double negation rule to get:

$$(\neg A \wedge B \wedge \neg B) \vee (A \wedge \neg B \wedge \neg B) \vee (A \wedge B \wedge B)$$

Applying the complement rule to the first bracket and the idempotent rules to the other two we simplify to:

$$(\neg A \wedge F) \vee (A \wedge \neg B) \vee (A \wedge B)$$

The first term is always false so we have

$$F \vee (A \wedge \neg B) \vee (A \wedge B) \equiv (A \wedge \neg B) \vee (A \wedge B)$$

by the identity rule. Now by the distributive rule we have $A \wedge (\neg B \vee B)$, and finally we can apply the complement and identity rules to get

$$A \wedge (\neg B \vee B) \equiv A \wedge T \equiv A.$$

The fact that only one clue is correct tells us that $A$ is true, that is the gold is in box A.

Why did we also need the second statement:

$$(A \wedge \neg B \wedge \neg C) \vee (\neg A \wedge B \wedge \neg C) \vee (\neg A \wedge \neg B \wedge C)?$$

Remember this tells us that there is gold in only one box. Since we know that the gold is in box A, it thus follows that there is no gold in boxes B and C. We don't really need the formula above to see this, but for completeness let's go ahead and use it.

Since $A$ is true and therefore $\neg A$ is false the statement becomes

$$(T \wedge \neg B \wedge \neg C) \vee (F \wedge B \wedge \neg C) \vee (F \wedge \neg B \wedge C)$$

The first term reduces to $\neg B \wedge \neg C$ by the identity rule, while the second and third reduce to $F$ giving us

$$(\neg B \wedge \neg C) \vee F \vee F \equiv \neg B \wedge \neg C$$

Hence (as expected) we conclude that $B, C$ are false i.e. the boxes B,C have no gold.

## 1.4   Implications

Consider the following truth-table where the truth value of the third column depends on the truth values of $P, Q$.

| $P$ | $Q$ | $*$ |
|-----|-----|-----|
| $T$ | $T$ | $T$ |
| $T$ | $F$ | $F$ |
| $F$ | $T$ | $T$ |
| $F$ | $F$ | $T$ |

Suppose we know that $*$ is true. What does this tell us about the relationship between $P$ and $Q$?

The $*$ statement is true if

*the statement $Q$ must be true whenever $P$ is true*

Thus this truth table encodes the notion of *implication:* $P \implies Q$.

We define the *implication* $P$ implies $Q$ or in symbols $P \implies Q$ to be the statement with truth table

| $P$ | $Q$ | $P \implies Q$ |
|---|---|---|
| $T$ | $T$ | $T$ |
| $T$ | $F$ | $F$ |
| $F$ | $T$ | $T$ |
| $F$ | $F$ | $T$ |

This can be expressed in terms of the fundamental logical connectives previously introduced by defining $P \implies Q$ as the statement $(\neg P) \vee Q$. Note that $(\neg P) \vee Q$ has exactly the truth table given above. Hence any statement involving $P \implies Q$ can be expressed as a disjunction using the logical equivalence

$$P \implies Q \equiv (\neg P) \vee Q.$$

By convention if we write something like $P \vee Q \implies R \wedge S$ we will interpret this as meaning $(P \vee Q) \implies (R \wedge S)$. In words we might say "if P is true or Q is true then R and S are both true": we can think of the *if...then...* as providing the implied brackets around then $P$ or $Q$.

**Example 1.6.** Consider the statement

> if $n - 1$ *is a multiple of* $4$ *then* $n^2 - 1$ *is a multiple of* $4$

We write this symbolically as $P(n) \implies Q(n)$ where:

$$P(n) : n - 1 \text{ is a multiple of } 4$$
$$Q(n) : n^2 - 1 \text{ is a multiple of } 4$$

The truth values of $P$ and $Q$ depend on the value of $n$

| $n$ | $P(n)$ | $Q(n)$ | $P(n) \implies Q(n)$ |
|---|---|---|---|
| $... - 3, 1, 5, ...$ | $T$ | $T$ | $T$ |
| | $T$ | $F$ | $F$ |
| $... - 1, 3, 7, ...$ | $F$ | $T$ | $T$ |
| $... - 2, 0, 2, 4, ...$ | $F$ | $F$ | $T$ |

For some values of $n$ both $P$ and $Q$ are true. For some values $P$ is false, but $Q$ is true. For some values (even $n$) both $P$ and $Q$ are false. But notice that for **all** values of $n$ the implication $P(n) \implies Q(n)$ is true, as there are

no values of $n$ for which $n - 1$ is a multiple of $4$ but where $n^2 - 1$ is not a multiple of $4$.

To reiterate, the meaning of $\implies$ is that whenever the first statement is true, the second statement must also be true.

**Example 1.7.** Consider the statement

*if $p$ is prime then $2^p - 1$ is prime*

How do we show that this is *false*?

Let us again consider a truth table:

| $p$ | $p$ is prime | $2^p$ is prime | $(p$ is prime$) \implies (2^p - 1$ is prime$)$ |
|---|---|---|---|
| $2, 3, 5, 7, 13, ...$ | $T$ | $T$ | $T$ |
| $11, 23, ...$ | $T$ | $F$ | $F$ |
| | $F$ | $T$ | $T$ |
| $1, 4, 6, 8, 9, ...$ | $F$ | $F$ | $T$ |

We disprove the statement by giving a *counterexample*, here $p = 11$ is a possible counterexample.

Counterexamples are cases where the **first statement is true** and the **second statement is false** which makes the implication **false**.

Symbolically we negate $P \implies Q$ as follows:

$$\neg(P \implies Q) \equiv \neg(\neg P \lor Q) \qquad \text{by definition of } \implies$$
$$\equiv (\neg\neg P) \land \neg Q \text{by de Morgan's rule}$$
$$\equiv P \land \neg Q \qquad \text{by the double } \neg \text{ rule}$$

The implication $P \implies Q$ is false when $P$ is true but $Q$ fails to be true.

## Converse and Contrapositive

The *converse* of an implication $P \implies Q$ is the reversed implication $Q \implies P$.

This has a different meaning to $P \implies Q$ and may be false even if $P \implies Q$ is true.

**Example 1.8.** $n - 1$ *is a multiple of* $4 \implies n^2 - 1$ *is a multiple of* $4$

is a true statement. Its converse is

$n^2 - 1$ *is a multiple of* $4 \implies n - 1$ *is a multiple of* $4$

which is false, for example $3^2 - 1$ is a multiple of $4$ but $3 - 1$ is not a multiple of $4$.

We define the *equivalence* $P \iff Q$ (or in words $P$ if and only if $Q$) to be the statement $(P \implies Q) \wedge (Q \implies P)$. This has truth table

| $P$ | $Q$ | $P \iff Q$ |
|:---:|:---:|:---:|
| $T$ | $T$ | $T$ |
| $T$ | $F$ | $F$ |
| $F$ | $T$ | $F$ |
| $F$ | $F$ | $T$ |

**Example 1.9.** $n^2 - 1$ *is a multiple of* $4 \iff n$ *is odd*

is a true statement.

The *contrapositive* of an implication $P \implies Q$ is the implication $\neg Q \implies \neg P$.

The implication $P \implies Q$ and its contrapositive $\neg Q \implies \neg P$ are equivalent: the first says

> *when $P$ is true then $Q$ must also be true*

but a consequence of this is

> *$Q$ is false could happen only if $P$ is false*

which is the contrapositive.

We can show that $P \implies Q \equiv \neg Q \implies \neg P$ as follows:

$$
\begin{aligned}
P \implies Q &\equiv \neg P \vee Q & \text{definition of } \implies \\
&\equiv Q \vee \neg P & \text{commutativity} \\
&\equiv \neg\neg Q \vee \neg P & \text{double } \neg \text{ rule} \\
&\equiv \neg Q \implies \neg P & \text{definition of } \implies
\end{aligned}
$$

**Example 1.10.** The statement

> $n^2 - 1$ *is a multiple of* $4 \implies n$ *is odd*

has contrapositive

> $n$ *is even* $\implies n^2 - 1$ *is not a multiple of* $4$

It is easy to see that the latter must be true since if $n$ is even then $n^2$ will be a multiple of $4$ and hence $n^2 - 1$ is not a multiple of $4$. Since the original statement is equivalent to its contrapositive, it follows that that statement is also true.

**Example 1.11.** Here's another word problem example.

*Andy, Billie, Charlie and Dannie are friends, some of whom are mathematicians. Consider the statements:*

1. *If Billie is not a mathematician then Andy and Charlie are not mathematicians*
2. *If Andy not a mathematician or Billie is a mathematician then Charlie is a mathematician.*
3. *If Billie and Charlie are mathematicians then neither Andy nor Dannie are mathematicians.*

*Assuming the statements are true, which of the friends are mathematicians?*

We begin by writing the statements symbolically:

1. $\neg B \implies (\neg A \wedge \neg C)$
2. $(\neg A \vee B) \implies C$
3. $B \wedge C \implies \neg(A \vee D)$

By definition of $\implies$ and the double $\neg$ rule, the first statement can be rewritten

$$\neg B \implies (\neg A \wedge \neg C) \equiv \neg\neg B \vee (\neg A \wedge \neg C) \equiv B \vee (\neg A \wedge \neg C).$$

Then by de Morgan's rule this is equivalent to

$$(B \vee \neg A) \wedge (B \vee \neg C).$$

As the statement is assumed to be true, both parts of the conjunction must be true, in particular $B \vee \neg A$ is a true statement. This is the same as $\neg A \vee B$ and by the second statement above we know $(\neg A \vee B) \implies C$. Hence we deduce that $C$ must be true.

But we also know that $B \vee \neg C$ is true so the conjunction of this with $C$ (that is $(B \vee \neg C) \wedge C$) is also true. So we have

$$(B \vee \neg C) \wedge C \equiv (B \wedge C) \vee (\neg C \wedge C) \equiv (B \wedge C) \vee F \equiv B \wedge C$$

applying the de Morgan, complement and identity rules.

We have therefore shown that $B$ and $C$ are true, but now from statement 3 it follows that the statement $\neg(A \vee D) \equiv \neg A \wedge \neg D$ holds.

Therefore Billie and Charlie are mathematicians while Andy and Dannie are not.

In the above example, a couple of times we used the observation that when we know that an implication $P \implies Q$ is true and also that $P$ itself is true, then we can conclude that $Q$ must also be true. (In the above example we did this using the statements $(\neg A \lor B) \implies C)$ and $\neg A \lor B$ to deduce $C$.). This principle goes by the name r defining("modus ponens") which can be justified by the following logical equivalences:

$$\begin{aligned} P \land (P \implies Q) &\equiv P \land (\neg P \lor Q) & \text{definition of } \implies \\ &\equiv (P \land \neg P) \lor (P \land Q) & \text{de Morgan's rule} \\ &\equiv F \lor (P \land Q) & \text{complement rule} \\ &\equiv P \land Q & \text{identity rule} \end{aligned}$$

These equivalences tell us that knowing $P$ holds and $P \implies Q$ is the same as knowing that $P$ and $Q$ are both true.

The following theorem gathers together some logical equivalences involving $\implies$ .

> **Theorem 1.12**
>
> - $P \implies Q \equiv \neg P \lor Q$ (definition)
> - $P \implies Q \equiv \neg Q \implies \neg P$ (contrapositive)
> - $P \land (P \implies Q) \equiv P \land Q$ (modus ponens)
> - $P \implies (Q \land R) \equiv (P \implies Q) \land (P \implies R)$
> - $(P \lor Q) \implies R \equiv (P \implies R) \land (Q \implies R)$
> - $(P \land Q) \implies R \equiv P \implies (Q \implies R)$

## 1.5 Proofs

What does it mean to prove a statement, and how do we go about this?

A *proof* is a sequence of steps used to demonstrate that a statement is true.

A proof will begin with a *premise*, meaning a collection of statements, or predicates, that we accept as true. This provides the context in which we are proving the statement. For example suppose we wanted to prove the statement:

*For an integer $n > 2$ there are no positive integer solutions $a, b, c$ to the equation*

$$a^n + b^n = c^n$$

Then our premise would be that $n$ is an integer greater than $2$ and that $a, b, c$ are positive integers.

We would then try to prove that $a, b, c$ do not give solutions to the equation $a^n + b^n = c^n$.

The result that we wish to show is true (when the premise is true) is called the *conclusion*.

A proof consists of a sequence of predicates each of which is true, under the assumptions of the premise. However we require more than just the condition that the predicates are true. Consider the following:

> *Let $n > 2$ be an integer and let $a, b, c$ be positive integers.*
>
> *Then $a^n + b^n \neq c^n$.*

Though the second statement is true, the above does not provide a proof of Fermat's Last Theorem because we have not justified *why* the second statement follows from the first.

For a sequence of predicates to be accepted as a proof, each must be:

- Something which we assume to be true: *a premise, definition, or axiom.*
- Something we already know to be true: *a statement we have previously proved.*
- An *additional assumption* meaning that we are considering a particular case. We are then beginning a *subproof* which will prove *part* of the original assertion: to complete the proof we must make sure that all possible cases are considered.
- A predicate which *follows directly* from one or more of the earlier steps of the proof.

We will now look at various examples of proofs. There may be steps in some proofs for which the mathematical ideas need further discussion at a later point, but for the moment our focus is on the shape and structure of proofs.

We also remark that the proofs in this section may seem rather long-winded. This is because we want to be absolutely precise in terms of what is and is not allowed in a proof. In practice one often omits some of the simple and routine steps from a proof, or combines steps, which is justified by the belief that one could "fill in the gaps" if challenged. At least for now, we will justify each logical step of a proof (though we may still abridge some algebraic manipulations).

## 1.5.1 Direct Proof

**Example 1.13.** Let us prove the statement

> *if $n$ is even then $n^2$ is even.*

Before we write a proof let's analyse the statement.

> The premise is "$n$ is even."
> The conclusion is "$n^2$ is even"

**Working backwards:**

Looking at the conclusion we ask the following key question:

> *How do we prove that a number is even?*

When asking this question we want to keep it general rather than asking "how do we prove the $n^2$ is even?" as this helps us to focus on the idea rather than the details.

To prove that a number is even we need to show that it can be written as $2$ times some other integer. So our aim is to find this integer.

Now that we have answered the general question we give a specific target statement. We aim to prove the assertion:

$$\text{there is an integer } l \text{ such that } n^2 = 2l$$

**Working forwards:**

Now looking at the premise, what do we know? Since $n$ is even, again it must be $2$ times some integer. We will give this integer the name $k$:

$$\text{there is an integer } k \text{ such that } n = 2k$$

When we write this statement we must choose a different name $k$ for this integer than we had in our target statement.

Note there is an important difference between the predicates we were considering: "$n$ *is even*" and "$n^2$ *is even*". The first is the premise so we **are allowed to take this as a true assertion**. The second is something that we are trying to prove, so **we cannot write it down as part of our proof** until we know that it is true.

What we *can* write down if we wish, since this is a true statement regarding the conclusion is the conditional:

> *if there is an integer $l$ such that $n^2 = 2l$ then $n^2$ is even*

or even

> *$n^2$ is even if and only if there is an integer $l$ such that $n^2 = 2l$.*

**Writing the proof:**

Here is the proof written out, indicating what the justification is for each step.

> $n$ is even *(premise)*
>
> there is an integer $k$ such that $n = 2k$ *(definition of even)*
>
> so $n^2 = (2k)^2$ *(substitution)*
>
> $n^2 = (2k)^2 = 4(k^2) = 2(2k^2)$ *(algebraic manipulation)*
>
> let $l = 2k^2$ *(define a variable called l)*
>
> then $n^2 = 2l$ *(substitution)*
>
> there is an integer $l$ such that $n^2 = 2l$ *(by construction)*
>
> $n^2$ is even *(definition of even)*

We have now proved the statement

> *if $n$ is even then $n^2$ is even.*

We will come back to the question of the allowed steps of algebraic manipulation later.

What about the substitution step? This involves the rules of *equality* which we will consider in the next section.

**Recap of the strategy**

The basic strategy for proving an "if *premise* then *conclusion*" statement is to begin by asking the question "what condition is required to prove a general statement of the form *conclusion*". From this we obtain a target statement to work towards (e.g. "there is an integer $l$ such that $n^2 = 2l$"). We then write down the *premise* and any statements which immediately follow from that. These statements typically involving the definitions of any predicates in the premise (e.g. we translate "is even" into the existence of a suitable $k$). At this point we try to work from the premise towards our towards our target statement. Note that although the target statement is the first thing we think about here it should appear at (or near) the **end** of the proof, not the beginning.

## 1.5.2 Equality

What does $=$ mean? When we write $a = b$ we mean that $a, b$ refer to the same "thing." Remember that when we discussed variables we said that a variable does not have to be a number. It might also refer to a set or a vector or something else.

When we write $a = b$ we mean that whatever type of object $a$ and $b$ refer to they are indistinguishable from one another.

Since equality can apply across all sorts of different mathematical objects we will consider the properties of this separately from other algebraic properties which will vary depending upon he context. (For instance: $ab = ba$ is true for integers, real numbers or even complex numbers, but $AB = BA$ is false for matrices; if $b \neq 0$ then $ax = b$ is equivalent to $x = a/b$ for real or rational numbers but $a/b$ is not defined for integers. Algebra is therefore specific to the collection of values which are allowed for the variables.)

Here are some key properties of equality.

- for any $a$, the statement $a = a$ is true
- for any $a$ and $b$, the statement $a = b$ is true if and only if $b = a$ is also true
- for any $a, b$ and $c$, if $a = b$ and $b = c$ are true then $a = c$ is true

We say that equality is:

- *reflexive* $a = a$
- *symmetric* $a = b \iff b = a$
- *transitive* $a = b \ \wedge \ b = c \implies a = c$

Equals has one other important property which is *substitution*:

- if $a = b$ then for any true statement involving $a$ we may replace one or more occurrences of $a$ with $b$.

This is the rule which allows us to conclude from $n = 2k$ (and of course $n^2 = n^2$) that $n^2 = (2k)^2$.

### 1.5.3 Proof by cases

We will now consider another example of a proof.

**Example 1.14.** Prove that if $a$ is even or $b$ is even then $ab$ is even.

As in the previous example we need to prove that a number ($ab$) is even and so must find a way to write this as $2$ times some integer.

Our premise is a disjunction $a$ *is even or* $b$ *is even*. When we have a disjunction we will often split into cases.

In outline the proof will be:

- show that if $a$ is even then $ab$ is even
- show that if $b$ is even then $ab$ is even

- combine these to say that if $a$ is even or $b$ is even then $ab$ is even

The first two steps involve "subproofs," that is we have a proof of another statement embedded in our proof:

We make an additional assumption (in the first case that $a$ is even) and the steps following this assumption are statements which are true **in this new context** but may not have been true in general (it might have been the case that $a$ was odd and $b$ was even instead). The conclusion at the end of the subproof is that $ab$ *is even in this context*. But now we can write a statement which is true given our original premise, namely that "**if** $a$ *is even then* $ab$ *is even*."

**Writing the proof:**

Here's the proof: in this example we use indentation to indicate where we are in one of the subproofs.

> $a$ is even or $b$ is even *(premise)*
> Suppose $a$ is even *(assumption)*
> > then there exists $k$ such that $a = 2k$ *(definition of even)*
> > so $ab = (2k)b$ *(substitution)*
> > $ab = 2(kb)$ *(associative rule)*
> > therefore $ab$ is even *(definition of even)*
> Hence if $a$ is even then $ab$ is even *(conclusion of the first case)*
> Suppose $b$ is even *(assumption)*
> > then there exists $l$ such that $b = 2l$ *(definition of even)*
> > so $ab = a(2l)$ *(substitution)*
> > $ab = 2(al)$ *(associative and commutative rules)*
> > therefore $ab$ is even *(definition of even)*
> Hence if $b$ is even then $ab$ is even *(conclusion of the second case)*
> So if $a$ is even then $ab$ is even and if $b$ is even then $ab$ is even *(combining the two statements)*
> Therefore if $a$ is even or $b$ is even then $ab$ is even

Our aim was to prove "$(a$ *is even* $\vee$ $b$ *is even)* $\implies$ $ab$ *is even*" but what we actually argued by proving the two cases was "$(a$ *is even* $\implies$ $ab$ *is even)* $\wedge$ $(b$ *is even* $\implies$ $ab$ *is even)*." However these statements are equivalent by the following logical equivalence: (see Theorem 1.12)

$$(P \vee R) \implies R \equiv (P \implies R) \wedge (Q \implies R)$$

### 1.5.4   Proof by contrapositive

In the same way that we used a logical equivalence in the previous example, if we want to prove an implication $P \implies Q$ it is equivalent to prove the contrapositive $\neg Q \implies \neg P$.

**Example 1.15.** Suppose we want to prove the following statement for integer values of $n$:

>   *if $n^2$ is even then $n$ is even.*

It is equivalent to prove the statement

>   *if $n$ is not even then $n^2$ is not even.*

As we are talking about integers we might rephrase this as

>   *if $n$ is odd then $n^2$ is odd.*

*Remark: If we take $n$ is odd to mean precisely that $n$ is not even then this is immediate. However instead we would like to take the view that $n$ is odd means there is an integer $k$ such that $n = 2k + 1$. The claim that $n$ is not even is the same as $n$ is odd therefore needs some justification which we will come back to later.*

Our task is now to prove that

>   *if $n$ is odd then $n^2$ is odd.*

so now the *premise* is $n$ is odd while we are aiming for the *conclusion* that $n^2$ is odd.

**Writing the proof:**

>   $n$ is odd *(premise)*
>
>   there is an integer $k$ such that $n = 2k + 1$ *(definition of odd)*
>
>   $n^2 = (2k + 1)^2$ *(substitution)*
>
>   $n^2 = 4k^2 + 4k + 1 = 2(2k^2 + 2k) + 1$ *(algebra)*
>
>   so $n^2$ is odd *(definition of odd)*
>
>   we have shown that if $n$ is odd then $n^2$ is odd
>
>   so if $n^2$ is even then $n$ is even *(logical equivalence and even $\equiv$ not odd)*

## 1.5.5 Proof by contradiction

Another important proof technique is proof by contradiction. For example if you have seen a proof that $\sqrt{2}$ is irrational then this probably involved proof by contradiction.

As we only want to work with *integers* in this course let us rephrase the statement as

*there are no positive integers $a, b$ satisfying $a^2 = 2b^2$*

(The equivalence of these two statements is justified by the fact that $\sqrt{2}$ is defined as the positive solution of the equation $x^2 = 2$ and so if $x = a/b$ then $a^2 = 2b$.)

The strategy of proof by contradiction is to assume that the statement $P$ to be proved is *false* and to derive from this a contradiction: a statement which must logically by false such as $Q \wedge \neg Q$).

The justification for this is that if we prove $\neg P$ (the statement that $P$ is false) implies a false statement then we have shown that the statement

$$\neg P \implies F$$

is true. But using logical equivalences

$$\neg P \implies F \equiv (\neg\neg P) \vee F \equiv P \vee F \equiv P$$

where the equivalences are respectively the definition of $\implies$, the double $\neg$ rule and the identity rule.

Hence proving $\neg P \implies F$ is equivalent to proving $P$.

**Example 1.16.** Prove that there are no positive integers $a, b$ satisfying $a^2 = 2b^2$.

**Writing the proof:**

Suppose $a, b$ are integers such that $a^2 = 2b^2$ *(premise)*

Let $d$ be the greatest common divisor of $a$ and $b$ *(define a variable $d$)*

$d$ is a divisor of $a$ and is a divisor of $b$ *(from the definition of greatest common divisor)*

There exist integers $u, v$ such that $a = du$ and $b = dv$ *(definition of divisor)*

$(du)^2 = 2(dv)^2$ *(substitution)*

$d^2 u^2 = d^2(2v^2)$ *(associativity and commutativity)*

$u^2 = 2v^2$ *(cancellation)*

$u^2$ is even *(definition of even)*

$u^2$ is even $\implies$ $u$ is even *(previous example)*

$u$ is even *(modus ponens)*

there exists $k$ such that $u = 2k$ *(definition of even)*

$(2k)^2 = 2v^2$ *(substitution)*

$2(2k^2) = 2v^2$ *(associative and commutative rules)*

$2k^2 = v^2$ *(cancellation)*

$v^2$ is even *(definition of even)*

$v$ is even *(previous example)*

there exists $l$ such that $v = 2l$ *(definition of even)*

$a = d(2k) = (d \cdot 2)k$ *(substitution and associativity)*

$b = d(2l) = (d \cdot 2)l$ *(substitution and associativity)*

$d \cdot 2$ is a divisor of $a$ *(definition of divisor)*

$d \cdot 2$ is a divisor of $b$ *(definition of divisor)*

$d \cdot 2$ is a common divisor of $a$ and $b$ *(definition of common divisor)*

$d = d \cdot 1 < d \cdot 2$ *(properties of $<$)*

$d$ is not the **greatest** common divisor of $a, b$ *(definition of greatest common divisor)*

$d$ is the greatest common divisor of $a, b$ *(definition of $d$)*

$d$ **is** the greatest common divisor of $a, b$ and $d$ is **not** the greatest common divisor of $a, b$ *(combining the previous two statements)*

Contradiction

Hence the initial assumption that $a, b$ are integers such that $a^2 = 2b^2$ is false.

The equation $a^2 = 2b^2$ has no integer solutions.

We will consider divisors, greatest common divisors, inequalities etc. in greater detail in a later chapter.

## 1.6   Sets and Quantifiers

### 1.6.1   Sets

Sets form an important part of the language of mathematics. Informally, a set simply means a collection of elements. These elements can be any type of mathematical object: they can be numbers, vectors, or even other sets.

Sets are denoted by curly brackets $\{...\}$.

For the purpose of this course, the most important sets are the set of natural numbers:

$$\mathbb{N} = \{1, 2, 3, \dots\}$$

and the set of integers

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

Note some authors may have the convention that the natural numbers start at $0$, however if we want to include $0$ then we will write

$$\mathbb{N}_0 = \{0, 1, 2, \dots\}$$

The notation $x \in X$ means that $x$ is an element of the set $X$. For example $n \in \mathbb{Z}$ means $n$ is an element of the set of integers or more briefly $n$ is an integer.

We can build *subsets* of a given set using *set builder notation*:

$$\{x \in X : P(x)\}$$

denotes the set of elements of $x$ such that the predicate $P(x)$ is true.  For example

$$\{n \in \mathbb{Z} : n \text{ is even}\}$$

would denote the set of even numbers. We read the colon : as *"such that,"* so we have $x \in X$ such that $P(x)$ is true, or $n \in \mathbb{Z}$ such that $n$ is even.

Note some authors use the notation $\{x \in X \mid P(x)\}$ instead of the colon.

If I write down a set $A = \{1, 2\}$ and another set $B = \{1, 2\}$ then, despite the different names, we think of these as the same set. This is encoded in the following rule (axiom) from set theory:

> *if all $x \in A$ are also in $B$ and all $x \in B$ are also in $A$ then $A$ and B are the same set*

This condition can also be expressed as saying that predicates $x \in A$ and $x \in B$ should be equivalent.

$$\text{if } x \in A \iff x \in B \text{ then } A = B.$$

If we just have the implication in one direction, for example $x \in A \implies x \in B$ then this mean that all elements of $A$ are also in $B$ but there might be elements of $B$ that are not in $A$.

If $x \in A \implies x \in B$ then we say that $A$ is a *subset* of $B$, written $A \subseteq B$ (or as $B \supseteq A$).

## 1.6.2 Quantifiers

Frequently we wish to consider statements that are true *for all* values of some given variable. For example

$$\text{for every integer } n, \ n^2 \geq 0.$$

In symbolic terms we write this as

$$\forall n \in \mathbb{Z}, n^2 \geq 0.$$

The symbol "$\forall$" is pronounced "for all" or "for every" or "for each", and is called the *universal quantifier*.

On the other hand, we may wish to consider statements that are only true *for some* values of some given variable. For example

$$\text{there exists (at least one) integer number } n \text{ such that } n^2 = 289.$$

In symbolic terms we write this as

$$\exists n \in \mathbb{Z}, n^2 = 289.$$

The symbol "$\exists$" is pronounced "there exists" or "for some", and is called the *existential quantifier*.

What would be the negation of the statement

$$\forall n \in \mathbb{Z}, n^2 \geq 0.$$

The statement says that "the square of every integer is non-negative" or "every integer has a non-negative square".

So the negation says that "not every integer has a non-negative square". In other words "there are some (at least one) integer(s) that do not have a non-negative square". In symbolic form this is

$$\exists n \in \mathbb{Z}, \neg(n^2 \geq 0).$$

In general, the negation of the universal quantifier

$$\forall x, P(x)$$

is the existential quantifier

$$\exists x, \neg P(x).$$

In the same way the negation of the existential quantifiers

$$\exists x, P(x)$$

is the universal quantifier

$$\forall x, \neg P(x).$$

**Example 1.17.** The negation of

$$\exists n \in \mathbb{Z}, n^2 = 289.$$

is

$$\forall n \in \mathbb{Z}, n^2 \neq 289.$$

This (false) statement asserts that all integers fails to be solutions of this equation.

## 1.6.3 Implicit quantifiers

Consider the (now familiar) statement

*if $n$ is even then $n^2$ is even*

In symbols we write this as the implication

$$n \text{ is even} \implies n^2 \text{ is even}$$

When we write this we are not referring to some specific $n$ for which this is true. Rather we are saying that the implication is **always** true.

Picking a specific value for $n$ we can ask about the truth values in this statement, for example if $n = 2$ then the statements $n$ is even and $n^2$ is even are both true, so the implication evaluates as $T \implies T$ which is true.

On the other hand picking another integer value for $n$, say $n = 3$ then we see that the statements $n$ is even and $n^2$ is even are both false. The implication becomes $F \implies F$ which is again true (by definition of $\implies$). Indeed this will work for any integer $n$: we know this because we proved that if $n$ is even then $n^2$ is even.

The original statement should perhaps more precisely be stated as

> *if $n$ is **any** even integer then $n^2$ is even*

which makes it clearer that the statement should be interpreted as

$$\forall n \in \mathbb{Z}, \ n \text{ is even} \implies n^2 \text{ is even}$$

**A note on brackets:** by convention the statement $\forall n, P(n) \implies Q(n)$ should be read as $\forall n, (P(n) \implies Q(n))$ rather than $(\forall n, P(n)) \implies Q(n)$.

Note that in our discussion of proofs we have looked at how to prove that a conclusion is **always true whenever the premise is true**. In other words we have really been proving statements with (implicit) universal quantifiers:

$$\forall n, P(n) \implies Q(n).$$

We did however also make use of existential quantifiers: every time we use the statement $n$ is even we expressed this as "there is an integer $k$ such that $n = 2k$" or in symbols:

$$\exists k \in Z, \ n = 2k$$

If we want to write the statement "if $n$ is even then $n^2$ is even" purely in symbols then we can do so as follows:

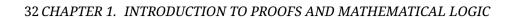$$\forall n \in \mathbb{Z}, (\exists k \in \mathbb{Z}, n = 2k) \implies (\exists l \in \mathbb{Z}, n^2 = 2l)$$

As per the above convention about brackets, the initial $\forall n \in \mathbb{Z}$ applies to the whole of this statement.

This proof that we gave for this statement demonstrates the use and proof of existential quantifiers. To prove an existential statement, which says there is *at least one value of the variable with a certain property*, we simply need to construct such an example.

In the case of our current statement, to prove $(\exists l \in \mathbb{Z}, n^2 = 2l)$ we need only to construct the value of $l$. This value is $2k^2$ where $k$ is the value that we know exists from the premise $(\exists k \in \mathbb{Z}, n = 2k)$.

We repeat here our earlier proof, but using quantifier notation:

> $n$ is even *(premise)*
> $\exists k \in \mathbb{Z}, n = 2k$ *(definition of even)*
> so $n^2 = (2k)^2$ *(substitution)*
> $n^2 = (2k)^2 = 4(k^2) = 2(2k^2)$ *(algebraic manipulation)*
> let $l = 2k^2$ *(define a variable called $l$)*
> then $n^2 = 2l$ *(substitution)*
> so $\exists l \in \mathbb{Z}, n^2 = 2l$ *(by construction)*
> $n^2$ is even *(definition of even)*

# Chapter 2

# The Integers

In this chapter we will begin our study of the **integers**.

## 2.1 Axioms

To prove results about integers we need a starting point. What are the integers? What properties do they have?

We must have some properties of the integers that we accept *without requiring a proof.* Without such a starting point, we would need to prove each property in terms of some previous result which in turn would require proof in terms of other results – and we would be stuck in an infinite descent.

These properties, which we accept without proof are referred to as *axioms*. We may think of these are defining what we mean by the set of integers.

In Chapter 3, we shall learn more of Euclid, who's most important contribution to mathematics, whilst also being famous for many other things, was the formulation of the axiomatic method.

In all of his writings he carefully sets out his assumptions and his rules for reasoning and then develops mathematical theory relying only on those rules or axioms.

We denote the set of integers by $\mathbb{Z}$. This comes from the German word *Zahlen*, meaning *numbers*.

The integers satisfy the following axioms:

- [Operations (Op)] The integers are equipped with two *binary operations*[1], addition and multiplication. These operations are denoted by the symbols $+, \cdot$ respectively.

- [Identity (Id)] $\mathbb{Z}$ contains two special elements, denoted $0$ and $1$, with $0 \neq 1$ such that for any integer $m$, $m+0 = 0+m = m$ and $m \cdot 1 = 1 \cdot m = m$.

- [Negation (Neg)] For every integer $m$, there exists an integer denoted $-m$ such that $m + (-m) = (-m) + m = 0$. We write $a - b$ as a shorthand for $a + (-b)$.

- [Commutative (Comm)] Addition and multiplication are commutative, i.e. $m + n = n + m$ and $m \cdot n = n \cdot m$ for all $m, n \in \mathbb{Z}$

- [Associative (As)] Addition and multiplication are associative, i.e. $(m + n) + p = m + (n + p)$ and $m \cdot (n \cdot p) = (m \cdot n) \cdot p$ for all $m, n, p \in \mathbb{Z}$

- [Distributive (Dist)] Multiplication distributes over addition, i.e, $m \cdot (n + p) = m \cdot n + m \cdot p$ for all $m, n, p \in \mathbb{Z}$.

- [Zero divisors (ZD)] If $m \cdot n = 0$ then $m = 0$ or $n = 0$. As well as the algebraic structure the integers have an ordering, by size, which interacts with the algebra. The axioms for this ordering are as follows:

- [Order relation (Ord)] $\mathbb{Z}$ is equipped with a relation $\leq$ called "less than or equal to".

- [Reflexive (Ref)] Every integer $m$ satisfies $m \leq m$.

- [Antisymmetric (ASy)] Given two integers $m, n$, if $m \leq n$ and $n \leq m$ then $m = n$.

- [Transitive (Tr)] If $m, n, p$ are integers with $m \leq n$ and $n \leq p$ then $m \leq p$.

- [Comparability (Comp)] Given any two integers $m, n$ either $m \leq n$ or $n \leq m$.

- [Shift (Sh)] If $m, n, p$ are integers with $m \leq n$, then $m + p \leq n + p$

- [Scale (Sc) ] If $m, n, p$ are integers with $m \leq n$ and $0 \leq p$ then $m \cdot p \leq n \cdot p$.

- [Well ordering (WO)] The *well ordering principle* – any non-empty subset of non-negative[2] integers has a unique least element.

We say that $n$ is the *least element* of a set $S$ if

- $n \in S$ and

---

[1]A binary operation on a set (here the integers $\mathbb{Z}$) means a function taking two elements from the set at input and returning one as output, for example $+ \; : \; (m, n) \mapsto m + n$.

[2]An integer $n$ is defined to be non-negative if $0 \leq n$.

- $\forall m \in S$ if $m \leq n$ then $m = n$.

The second condition is really telling us that **no $m \in S$ is strictly less than $n$.**

We accept these axioms as being true without requiring any justification of them. We will however justify almost **ALL** other results in this course.

We can also introduce the symbols $\geq, <$ and $>$ defined as:

- $m \geq n$ is defined to be $n \leq m$
- $m < n$ is defined to be $m \leq n$ *and* $m \neq n$
- $m > n$ is defined to be $n \leq m$ *and* $m \neq n$

Note that if $m < n$ is *false* then $m \leq n$ is false or $m \neq n$ is false. If $m \leq n$ is false then by *(Comp)* we have $n \leq m$. If $m \neq n$ is false then (by double $\neg$) $m = n$ is true so $n \leq m$ by (Ref). Hence *not $m < n$* implies $n \leq m$.

Conversely suppose $n \leq m$ is *true.* The predicate $m \leq n$ is either true or false. If it is true then $m = n$ by *(Ref)* so $m < n$ is false. On the other hand if $m \leq n$ is false then $m < n$ is false by definition.

Hence we have proved that "*not $m < n$*" $\iff$ "$n \leq m$". Of course it also follows that "*not $n \leq m$*" $\iff$ "$m < n$". We can therefore freely use this "translation" of *not $n \leq m$* in any future proofs.

---

**Note 2.1**

None of the axioms mention division. Indeed we don't even have a symbol for division since this is not defined as an operation on pairs of integers.

What *can* we do if we want to divide integers? We can take the *integer part* of the quotient along with the *remainder*, for example $13$ divided by $5$ is $2$ with remainder $3$. We do not need the notation of division to write this, instead we can write

$$13 = 2 \cdot 5 + 3.$$

The fact that we can always "divide" (be a non-zero integer) to get a quotient and remainder is a result called the **Remainder Theorem** which we will prove at the end of this chapter.

---

**Example 2.2.** Show that if $x$ is a non-zero integer, then the sequence

$$x, 2x, 3x, \ldots$$

contains no repetition.

The example is asking us to show that if $a$ and $b$ are distinct positive integers then $ax \neq bx$. (In fact we won't use positivity here as this is not required.)

The strategy is to use a proof by contradiction. Supposing that $ax = bx$ we would like to rearrange this as $(a - b)x = 0$ and then use *(ZD)* to deduce that either $a = b$ or $x = 0$ both of which are assumed to be false ($x$ is non-zero and $a, b$ are distinct.)

How do we rearrange $ax = bx$ to get $(a - b)x = 0$? We would like to subtract $bx$ from both sides, or more precisely add $-(bx)$ to both sides. However this yields the left-hand side as $(ax) + -(bx)$ which – as written – does not have a common term of $x$ to pull out using *(Dist)*.

To fix this we should add $(-b)x$ to both sides[3]. This will leave us with $bx + (-b)x = (b + -b)x = 0 \cdot x$ on the right-hand side, and this is zero.

**But wait**, how do we know that $0 \cdot x = 0$? We need to prove this first. *Note: (ZD) tells us that if $m \cdot n = 0$ then $m = 0$ or $n = 0$ but we need the converse: if $m = 0$ or $n = 0$ then $m \cdot n = 0$.*

> **Lemma 2.3**
>
> [a] For all integers $n$
> $$0 \cdot n = n \cdot 0 = 0.$$
> _____
> [a]A Lemma is a little theorem which will be used to prove other results.

*Proof.* We will focus on showing that $0 \cdot n = 0$ since we can use *(Comm)* to show that $n \cdot 0 = 0 \cdot n$.

When we look at the axioms we see that there is only one axiom *(Id)* which allows us to write something which is a product as an expression which does not involve a product. We will therefore need to use this axiom:

$$1 \cdot n = n$$

The idea is now to subtract $n$ from both sides. We then need to show that $1 \cdot n + (-n)$ is the same as $0 \cdot n$ for which the trick is to write $1 \cdot n$ as $(0 + 1) \cdot n$

_____
[3]Of course it is true that $-(bx) = (-b)x$ but we would need to prove this first.

and use distributivity:

$$
\begin{aligned}
0 \;=\;& n + (-n) & \text{(Neg)} \\
=\;& (1 \cdot n) + (-n) & \text{(Id)} \\
=\;& ((0 + 1) \cdot n) + (-n) & \text{(Id)} \\
=\;& ((0 \cdot n) + (1 \cdot n)) + (-n) & \text{(Dist)} \\
=\;& (0 \cdot n) + ((1 \cdot n) + (-n)) & \text{(As)} \\
=\;& (0 \cdot n) + (n + (-n)) & \text{(Id)} \\
=\;& (0 \cdot n) + 0 & \text{(Neg)} \\
=\;& 0 \cdot n & \text{(Id)}
\end{aligned}
$$

$\square$

We can now return to the solution to Example 2.2.

*Solution.* Our premise is that $x \neq 0$ and we also suppose, by way of contradiction, that $ax = bx$.

Then adding $(-b)x$ to both sides we deduce that

$$ax + (-b)x = bx + (-b)x.$$

Now the right hand side of this equation is equal to $(b + (-b))x$ by *(Dist)* which equals $0 \cdot x$ by *(Neg)*. By Lemma 2.3 we know that $0 \cdot x = 0$.

Hence

$$ax + (-b)x = bx + (-b)x = 0.$$

Now rearranging the left-hand side $ax + (-b)x$, by axiom *(Dist)* becomes $(a + (-b))x$, so we have shown that $(a + (-b))x = 0$.

From axiom *(ZD)*, we can now deduce that either $a + (-b) = 0$ or $x = 0$. By our original assumption, $x \neq 0$ and so we conclude that $a + (-b) = 0$. We can now add $b$ to both sides of this equation to deduce

$$(a + (-b)) + b = 0 + b = b$$

by *(Id)*. Now by *(As)* the left-hand side is $a + ((-b) + b)$ and applying *(Neg)* we have

$$a + 0 = a + ((-b) + b) = b$$

Consequently we deduce that $a = b$ by axiom *(Id)*.

This is the contradiction that we are seeking as $a, b$ were assumed to be distinct and so we conclude that the terms in the sequence are all distinct.

When when we have an equation like $a + (-b) = 0$ we often wish to add $b$ (to get $a = b$). To save going through the axioms every time, let's write a Lemma we can come back to.

> **Lemma 2.4**
>
> For all integers $m, n$
> $$(m + (-n)) + n = m$$
> and
> $$(m + n) + (-n) = m$$

*Proof.* We have $(m + (-n)) + n = m = m + ((-n) + n)$ by *(As)*.

Now $(-n) + n = 0$ by *(Neg)* so $(m + (-n)) + n = m = m + 0$.

This equals $m$ by *(Id)*.

The second statement is proved in exactly the same way!     □

## 2.2  Decimal Expansions

When we think about numbers we usually think of a list of digits e.g. $2024$ or $1001$.

*How does this picture of the integers compare with our description in terms of the axioms?*

We would like to prove that every positive integer has a decimal expansion. We will give a proof of this later. Let us begin by asking what the "decimal expansion" of an integer means.

The axioms introduce the numbers $0$ and $1$, and we define the other possible digits $2, 3, \ldots, 9$ by

| $2 :=$ | $3 :=$ | $4 :=$ | $5 :=$ | $6 :=$ | $7 :=$ | $8 :=$ | $9 :=$ |
|--------|--------|--------|--------|--------|--------|--------|--------|
| $1 + 1$ | $2 + 1$ | $3 + 1$ | $4 + 1$ | $5 + 1$ | $6 + 1$ | $7 + 1$ | $8 + 1$ |

When we write $2024$ this is a list of digits $2, 0, 2, 4$ which in reverse order are the units, tens, hundreds, thousands (and so on if we have a larger number). So really we are writing the number as a sum $2 \cdot 10^3 + 0 \cdot 10^2 + 2 \cdot 10^1 + 4 \cdot 10^0$ (where of course $10 := 9 + 1$):

| Thousands | Hundreds | Tens | Units |
|-----------|----------|------|-------|
| 2 | 0 | 2 | 4 |
| $2 \cdot 10^3$ | $0 \cdot 10^2$ | $2 \cdot 10^1$ | $4 \cdot 10^0$ |

For an integer $n$ its *decimal expansion* is a list of numbers $d_k, d_{k-1}, \ldots, d_0$ with the properties that + each digit $d_i$ is an integer such that $0 \le d_i \le 9$ + the first digit $d_k$ is non-zero + $n$ is equal to the following sum:

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

The notation $\sum_{i=0}^{k} d_i \cdot 10^i$ means the sum $d_0 \cdot 10^0 + d_1 \cdot 10^1 + \cdots + d_{k-1} \cdot 10^{k-1} + d_k \cdot 10^k$. Of course we have also introduced the notation of powers to indicate repeated multiplication.

Note the following standard convention on powers:

$$a^{b^c} \text{ means } a^{(b^c)}$$

The reason why we interpret $a^{b^c}$ in this way instead of as $(a^b)^c$ is that $(a^b)^c$ can be rewritten without brackets as $a^{bc}$.

Implicit in the definition of the decimal expansion for a positive integer is the fact that the digits $1, 2, \ldots 9$ are all positive. How do we know this? We'll show that $0 < 1$, from which we could also show that $0 < 2, 3$ etc.

## 2.3 Order properties of $\mathbb{Z}$

When we say that a number $n$ is *positive* we mean that it satisfies the strict inequality $0 < n$; in other words $0 \le n$ and $0 \ne n$.

The set $\mathbb{N}$ of natural numbers is the set of all integers $n$ such that $0 < n$.

---

**Lemma 2.5**

The number $0$ is in $\mathbb{N}$, that is $0 < 1$.

---

We will use proof by contradiction. Recall that we showed that $\neg(0 < 1)$ is $1 \le 0$.

*Proof.* Suppose $1 \le 0$. Then $1 + (-1) \le 0 + (-1)$ by *(Sc)*.

So $0 \le -1$ by *(Neg),(Id)*.

As $1 \le 0$ and $0 \le -1$ we have

$$1 \cdot (-1) \le 0 \cdot (-1)$$

by *(Sc)*. Hence $-1 \le 0$ by *(Id)* and Lemma 2.3.

But we also showed $0 \leq -1$ so $0 = -1$ by *(ASy)*.

Hence $1 = 0 + 1 = -1 + 1 = 0$ by *(Id),(Neg)*.

This contradicts $0 \neq 1$ which is true by *(Id)*.

We deduce that $1 \leq 0$ is false, so $0 < 1$ is true. □

We will show now show that $1$ is the *least positive integer*.

There must be a least positive integer by the axiom *(WO)*. Recall this axiom says:

> The *well ordering principle* – any non-empty subset of non-negative integers has a unique least element.

---

**Theorem 2.6**

$1$ is the least element of the set $\mathbb{N}$ of positive integers.

---

*Proof.* Let $n$ be the least element of $\mathbb{N}$, which exists by $(WO)$.

Our strategy is to show that $n$ cannot be greater than $1$ or less than $1$. By *(Comp)* we know that $1 \leq n$ or $n \leq 1$. We will show that in each case we deduce that $n = 1$.

Suppose $1 \leq n$. Since $1 \in \mathbb{N}$ and $n$ is the *least element* of $\mathbb{N}$ we have $1 = n$ as required.

Suppose $n \leq 1$. As $n \in \mathbb{N}$ we have $0 \leq n$ so by *(Sc)*

$$n \cdot n \leq 1 \cdot n = n.$$

Now $0 \leq n \cdot n$ by *(Sc)* since $0 \leq n$.

If $0 = n \cdot n$ then $n = 0$ or $n = 0$ by *(ZD)* so $n = 0$ (this is the idempotent rule from logic). But $n$ is in $\mathbb{N}$ so $n \neq 0$. This is a contradiction so we deduce that $0 < n \cdot n$.

Thus $n \cdot n \in \mathbb{N}$ and $n \cdot n \leq n$ so as $n$ is the *least* element of $\mathbb{N}$ we have $n \cdot n = n$.

Rearranging this $n \cdot n = n = 1 \cdot n$ by *(Id)* so

$$(n - 1) \cdot n = n \cdot n + (-1) \cdot n = 1 \cdot n + (-1) \cdot n = 0 \cdot n$$

by *(Dist)* twice. This is zero by Lemma 2.3 so by *(ZD)* either $n - 1 = 0$ or $n = 0$. As before we can rule out the case $n = 0$ so $n - 1 = 0$, and rearranging

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

using Lemma 2.4, *(Id)*.

So we have proved that if $1 \leq n$ then $n = 1$ and if $n \leq 1$ then $n = 1$. Hence we deduce that the least element $n$ of $\mathbb{N}$ is $1$. □

---

**Corollary 2.7**

$$\forall m, n \in \mathbb{Z}, \ m < n \iff m + 1 \leq n$$

---

*Proof.* Suppose that $m < n$. Then $m \leq n$ so $0 = m - m \leq n - m$ by *(Neg)*,*(Sh)*, and if $n - m = 0$ then $n = (n - m) + m = 0 + m = m$, by Lemma 2.4, *(Id)*.

So $n - m$ is strictly positive, therefore by the theorem $1 \leq n - m$. Hence

$$1 + m \leq (n - m) + m = n$$

by *(Sh)* and Lemma 2.4. Hence $m + 1 \leq n$ be *(Comm)*.

Conversely suppose $m + 1 \leq n$. Reversing the above argument:

$$1 = (1 + m) - m \leq n - m.$$

by Lemma 2.4, *(Sh)*.

As $0 \leq 1$ and $1 \leq n - m$ we have $0 \leq n - m$ by *(Tr)*. Thus $m = 0 + m \leq n - m + m = n$ by *(Id)*,*(Sh)* and our favourite lemma.

If $m = n$ then $m + 1 \leq m$ which implies $1 \leq 0$ by *(Sh)*. This contracts $0 < 1$, so $m \neq n$. Hence the inequality is strict $m < n$. □

# 2.4 Proof by Induction

We shall see a number of other examples of using the axioms to prove basic results in number theory later in the course. For now, we conclude by introducing a final method of proof, that many of you will have seen before.

This method is referred to as *Proof by Induction* and concerns proving statements of the form

$$\forall n \in \mathbb{N}, \ P(n).$$

The method of proof uses the Principle of Mathematical Induction, which can be expressed as

> **Theorem 2.8: Principle of Mathematical Induction**
>
> If $P(1) \wedge [\forall n \geq 1,\ P(n) \implies P(n+1)]$ then $P(n)$ is true for all $n \geq 1$.

So to prove a statement of the form

$$\forall n \geq 1,\ P(n).$$

we prove

1. $P(1)$ is true
2. $\forall n \geq 1,\ P(n) \implies P(n+1)$.

The first step is often referred to as the *base step* or *anchoring step*, while the second step is referred to as the *inductive step*. To prove the inductive step we take $P(n)$ as premise and prove $P(n+1)$. The premise $P(n)$ is referred to as the *induction hypothesis*.

We shall prove this theorem shortly but it is helpful to look at a simple example first.

The principle relies on the fact that if a statement about the positive integers is false for some positive integers, then there is a first one at which it fails to be true. If we take $n$ to be the last number before it fails then $P(n)$ is true but $P(n+1)$ is false, so if we want to prove that a statement is true for all positive integers we just have to show that when $P(n)$ is true then $P(n+1)$ is also true.

**Example 2.9.** The equation $\displaystyle\sum_{i=1}^{n} 2i = n \cdot (n+1)$ seems to be true for all values of $n$. How would we prove this?

$P(1)$, the statement that $2.1 = 1.2$, is true by *(Comm)*.

Now if $\displaystyle\sum_{i=1}^{n} 2i = n \cdot (n+1)$ fails for some values, let $n$ be one less than the first value where it fails. So $P(n)$ is true but $P(n+1)$ is false. Therefore $\displaystyle\sum_{i=1}^{n} 2i = n \cdot (n+1)$ and so

$$\sum_{i=1}^{n+1} 2i = 2(n+1) + \sum_{i=1}^{n} 2i = 2(n+1) + n \cdot (n+1) = (n+1)(n+2)$$

where the first inequality follows from the definition of $\sum$, the second as $P(n)$ is true, and the third is *(Dist)*.

But the equation $\sum_{i=1}^{n+1} 2i = (n+1)(n+2)$ tells us that $P(n+1)$ is true! This is a contradiction so there cannot be a first place at which the statement becomes false, and it must be true for all values of $n$.

Here is a proof of Theorem 2.8.

*Proof.* We will use proof by contradiction so our premise is

$$P(1) \wedge [\forall n \geq 1, \ P(n) \implies P(n+1)]$$

and we suppose (for a contradiction) that

$$\exists k \geq 1, \ P(k) \text{ is false}.$$

(We've used a different variable to avoid confusion with $n$ in the previous statement.)

As in 2.9 the idea is to consider $n$ to be one less than the first value where $P(k)$ is false.

By assumption $P(k)$ fails for some $k$ so

$$S = \{k \in \mathbb{Z} : k \geq 1 \text{ and } P(k) \text{ is false}\}$$

is a non-empty subset of $\mathbb{N}$. By *(WO)* this set has a least element $m$.

As $m \in S$ we know $1 \leq m$. We consider cases:

If $m = 1$ then $P(1)$ is false, contradicting the premise.

If $m \neq 1$ then $1 < m$ so $1 + 1 \leq m$ by Corollary 2.7.

Applying *(Sh)* and Lemma 2.4 we get $1 \leq m - 1$.

Let $n = m - 1$. Then $n < n + 1$ (by *(Sh)* as $0 < 1$).

But $n + 1 = m$ which is the least element of the set $S$.

So $n \geq 1$ but $n \notin S$ meaning that $P(n)$ is true.

But $P(n+1) = P(m)$ which is false as $m \in S$.

So we have $P(n) \wedge \neg P(n+1) \equiv \neg(P(n) \implies P(n+1))$.

We have shown $\exists n \in \mathbb{N}, \ \neg(P(n) \implies P(n+1))$, but this contradicts the premise as

$$\exists n \geq 1, \ \neg(P(n) \implies P(n+1)) \equiv \neg [\forall n \geq 1, \ P(n) \implies P(n+1)]$$

Since both cases ($m = 1$ and $m \neq 1$) lead to contradictions the assumption (that $P(k)$ fails for some $k$) is false, meaning

$$\forall k \in \mathbb{N}, \; P(k) \text{ is true} .$$

$\square$

Notice also that there is nothing special about the starting value of 1. If we wish to prove the statement $\forall n \geq m \; P(n)$ we do so by proving

1. $P(m)$ is true
2. $\forall n \geq m \; P(n) \implies P(n + 1)$. Here $m$ could be positive, zero, or even negative.

**Example 2.10.** Show that for all $n \geq 1$

$$n! \leq n^n$$

*Solution.* The base step is to prove the statement

$$P(1) : 1! \leq 1^1$$

This is quite obvious, as both sides are equal to 1.

The inductive step is to prove the statement

$$\forall n \geq 1 \; P(n) \implies P(n + 1),$$

or in other words

$$\forall n \geq 1 \; n! \leq n^n \implies (n + 1)! \leq (n + 1)^{n+1}.$$

The best strategy with this kind of example is to take one side of the inequality in the statement $P(n + 1)$ and rewrite it using the corresponding expression from the statement $P(n)$.

Assume $n! \leq n^n$ (induction hypothesis)

$$
\begin{aligned}
(n + 1)! & = & (n + 1)n! & \quad \text{by definition of } n! \\
& \leq & (n + 1)n^n & \quad \text{by (Sc) and the induction hypothesis} \\
& < & (n + 1)(n + 1)^n & \quad \text{by (Sc) as } n < n + 1 \\
& = & (n + 1)^{n+1} & \quad \text{by definition of powers}
\end{aligned}
$$

Hence the result follows by mathematical induction.

**Exercise 2.11.** Show that for all $n \geq 1$

$$1 + 3 + \ldots + (2n - 1) = n^2$$

**Exercise 2.12.** Show that for all $n \geq 1$

$$(1 + a)^n \geq 1 + na$$

where $a$ is any real number greater than -1.

# Strong Induction

There is a more general form of induction called *strong induction*. In the standard form of induction (sometimes called weak induction) we show that when a result holds for $n$ then it also holds for $n + 1$. For **strong** induction we show that if $P(m)$ is true **for all** $m < n + 1$ then $P(n + 1)$ is true.

The idea is that this should make it easier to prove $P(n + 1)$ as we can use **all of the statements** $P(1), P(2), \ldots, P(n)$.

Strong induction is really just induction where the the induction hypothesis is that $P(1), P(2), \ldots, P(n)$ are **all** true.

In the induction step we must therefore prove that $P(1), P(2), \ldots, P(n), P(n+1)$ are **all** true. But since $P(1), P(2), \ldots, P(n)$ are true by assumption we only need to show that $P(n + 1)$ is true.

In a strong induction we may start with multiple cases in the base step, i.e. we prove $P(1), P(2), \ldots, P(k)$ directly for some (hopefully small) $k$, and then start the induction at $k$.

**Example 2.13.** Let $f_n$ be the $n^{th}$ Fibonacci number[4] and let $\alpha = (1+\sqrt{5})/2^5$. Show that for $n \geq 3$

$$f_n > \alpha^{n-2}.$$

*Solution.* The base step in this case, involves proving the statement for $n = 3$ and $n = 4$. Since $\alpha < 2 = f_3$ and $\alpha^2 = (3 + \sqrt{5})/2 < 3 = f_4$, the result is true for $n = 3$ and $n = 4$.

Suppose then that $f_k > \alpha^{k-2}$ for $k = 3, 4, \ldots, n$ with $n \geq 4$ (this is the induction hypothesis) and consider the case $k = n + 1$.

We know in particular that $f_n > \alpha^{n-2}$ and $f_{n-1} > \alpha^{(n-1)-2} = \alpha^{n-3}$.

It is easy to check that $\alpha^2 = \alpha + 1$. Hence

$$\alpha^{(n+1)-2} = \alpha^{n-1} = \alpha^2 \alpha^{n-3} = (\alpha + 1)\alpha^{n-3} = \alpha^{n-2} + \alpha^{n-3}.$$

---

[4]The Fibonacci sequence is defined by $f_1 = f_2 = 1$ and $f_n = f_{n-1} + f_{n-2}$ for $n \geq 3$.
[5]This, of course, is a statement about real numbers not integers, but is included here as an example of strong induction.

So

$$f_{n+1} = f_n + f_{n-1} > \alpha^{n-2} + f_{n-1} > \alpha^{n-2} + \alpha^{n-3} = \alpha^{n-1}$$

using the inequalities for $f_n, f_{n-1}$ from the induction hypothesis.

We have thus shown that $f_k > \alpha^{k-2}$ for $k = 3, 4, \ldots, n + 1$, and hence by induction $f_k > \alpha^{k-2}$ for all $k \geq 3$.

## 2.5   The Remainder Theorem

As we noted at the beginning of this chapter, while we cannot divide integers, we can take the *integer part* of the quotient along with the remainder.

For example $13$ divided by $5$ is gives a quotient of $2$ with a remainder of $3$ which we can write as

$$13 = 2 \cdot 5 + 3.$$

We will now prove the following theorem:

---

**Theorem 2.14: Remainder Theorem**

If $a$ and $b$ are integers with $a \geq 0$, $b > 0$, then there is a **unique** pair of integers $q$ and $r$ such that $a = qb + r$ and $0 \leq r < b$.

---

NOTES:

- We refer to the value $q$ as the *quotient* and $r$ as the *remainder*.
- This result is also referred to as the *division algorithm*.
- Without the constraint $0 \leq r < b$ we could still assert *existence* of $q$ and $r$, but we could not assert *uniqueness*.
- What happens if we drop the constraints $a \geq 0$, $b > 0$?

*Proof.* Theorem 2.14 asserts that something exists and that it is unique. We have to prove both:

**Existence:** We will prove this by induction. Note we have two independent variables $a, b$, but we suppose that $b$ has been chosen and use induction on $a$. (Since $b$ can be chosen as any positive value, we will have proved this for **all** $b$ as well as **all** a).

Base step: For $a = 0$ there is s solution of $a = bq + r$ given by $q = r = 0$. (Using Lemma 2.3.)

Induction step: We assume the result holds for $a$, that is: there exist integers $q$ and $r$ such that $a = qb + r$ and $0 \leq r < b$.

Then $a + 1 = (qb + r) + 1 = qb + (r + 1)$.

$r < b$ so $r + 1 \leq b$ by Corollary 2.7.

We now divide into cases: $r + 1 = b$ and $r + 1 \neq b$ (that is $r + 1 < b$).

Suppose $r + 1 < b$. Then $a + 1 = qb + r'$ where $r' = r + 1$ and $r' < b$.

Note that $0 \leq r$ and $r \leq r + 1 = r'$ (since $0 < 1$) hence $0 \leq r'$ by *(Tr)*.

Hence if $r + 1 < b$ then we have shown that there exists $q, r' \in \mathbb{Z}$ such that $a + 1 = qb + r'$ and $0 \leq r' < b$, as required.

Now suppose $r + 1 = b$. Then

$$a + 1 = qb + b = qb + 1b = (q + 1)b$$

by *(Id),(Dist)*.

Using *(Id)* again $a + 1 = (q+1)b + 0$, so we have $a + 1 = q'b + r'$ where $q' = q + 1$, $r' = 0$ and $0 \leq r' < b$.

Again we have shown that for $a + 1$ we can find the quotient and remainder as required.

Hence by induction it follows that for **all** $a \geq 0$ there exist integers $q$ and $r$ such that $a = qb + r$ and $0 \leq r < b$.

**Uniqueness:** Suppose that

$$a = qb + r = q'b + r' \text{ with } 0 \leq r < b \text{ and } 0 \leq r' < b.$$

Then, applying Lemma 2.4, *(As)*, *(Comm)*,

$$qb - q'b = qb + r - r - q'b = q'b + r' - r - q'b = r' - r.$$

As $0 \leq r$ we have $-r \leq 0$ *(Sh)* so $r' - r \leq r' + 0 = r' < b$.

If $q - q' > 0$ then $q - q' \geq 1$ by Theorem 2.6 so

$$r' - r = qb - q'b = (q - q')b \geq 1b = b$$

by *(Dist),(Sc),(Id)*, which is a contradiction. Hence $q - q' \leq 0$ so $q \leq q'$ by *(Sh)*.

But likewise $q'b - qb = r - r'$ and $r - r' < b$. So if $q' - q > 0$ then we have a contradiction. We deduce that $q' \leq q$.

Since $q \leq q'$ and $q' \leq q$ we have $q = q'$ by *(ASy)*. Thus $q$ is unique.

Moreover

$$r - r' = (q' - q)b = 0b = 0$$

by Lemma 2.3 so $r = r - r' + r' = 0 + r' = r'$. Hence $r$ is also unique. $\qquad \square$

## 2.6   More Decimal Expansions

We finish this chapter by using the Remainder Theorem to prove that every positive integer has a decimal expansion.

There is nothing special about the number $10$ here, we could of course use any base $b > 1$.

---

**Theorem 2.15**

Every positive integer $n$ has a decimal expansion.

---

Note that the decimal expansion is unique. We leave this as an exercise.

*Proof.* We will prove this using strong induction on $n$.

Base cases: If $n < 10$ then $n$ is its own decimal expansion: $n = d_0 \cdot 10^0$ where $d_0 = n$ and $0 < d_0 \leq 9$.

Induction step: We take $n \geq 10$ and assume inductively that every integer from $1$ to $n - 1$ has a decimal expansion.

By the Remainder Theorem there exist integers $q, r$ such that

$$n = 10q + r \text{ and } 0 \leq r < 10$$

By 2.7 we can say $r + 1 \leq 10$ so $r \leq 9$.

As $n \geq 10$ we must have $q > 0$ so $q \geq 1$ by 2.6.

So $10q = (9 + 1)q = 9q + q \geq 9 + q > q$ by *(Sc),(Sh)* and using $0 < 1$.

As $0 \leq r$ we have $10q \leq n$ and we deduce that $q < n$.

Hence by the induction hypothesis $q$ has a decimal expansion:

$$q = d_k \cdot 10^k + \cdots + d_1 \cdot 10^1 + d_0 \cdot 10^0$$

where each digit $d_i$ is an integer such that $0 \leq d_i \leq 9$ and $d_k \neq 0$.

By *(Dist),(Comm)*:

$$10q = d_k \cdot 10^{k+1} + \cdots + d_1 \cdot 10^2 + d_0 \cdot 10^1$$

so $n = d_k \cdot 10^{k+1} + \cdots + d_1 \cdot 10^2 + d_0 \cdot 10^1 + r$.

Hence we have shown that there exists numbers $d'_0, \ldots d'_{k+1}$ defined by

$$d'_0 = r, \quad d'_{i+1} = d_i \text{ for } i = 0, \ldots, k$$

such that

$$n = d'_{k+1} \cdot 10^{k+1} + \cdots + d'_2 \cdot 10^2 + d'_1 \cdot 10^1 + d_0 \cdot 10^0$$

where each digit $d'_i$ is an integer such that $0 \leq d'_i \leq 9$ and $d'_{k+1} \neq 0$. (For $i > 0$ this holds because we know this for $d_{i-1}$ while for $d_0$ we showed that $0 \leq r \leq 9$.)

Hence we have shown that if all positive integers less than $n$ have decimal expansions then $n$ also has a decimal expansion.

By induction it follows that all positive integers have decimal expansions.

$\square$

# Chapter 3

# Divisibility and Euclid's Algorithm

## 3.1 Divisors

We all know that $2$ is a *divisor* or *factor*[1], of $6$, but that it is NOT a divisor of $7$. It is tempting to say that this is because $6/2$ is a whole number but $7/2$ is not, and if we are working in the world of fractions (or the real numbers or the complex numbers) this would make sense. However the very fact that $7/2$ is not an integer means that the operation of division is not defined in general for pairs of integers unless we do work in these worlds.

In this course we want to work inside the set of integers without reference to the larger number systems, so we won't usually be allowed to divide and we therefore try to avoid using arguments involving the concept of division, wherever possible. This raises the question of how to define a factor or a divisor without referring to division. We do this as follows:

> **Definition 3.1**
>
> We say that an integer $d$ is a *divisor* of the integer $a$ if and only if there is an integer $b$ such that $db = a$. We write $d|a$ as a shorthand for the statement '$d$ is a divisor of $a$', which we are also allowed to phrase as '$d$ divides $a$'.

---

[1]The words divisor and factor mean the same thing. We will tend to use the word divisor when talking about integers, though we may talk about factors in the context of algebra e.g. $x + 1$ is a factor of $x^2 - 1$.

**Example 3.2.** Since $2 \cdot 3 = 6$, $2$ is a divisor of $6$. On the other hand if we consider the equation $2b = 7$ we see that the left hand side is even and the right hand side is odd so there are no solutions, hence $2$ is not a divisor of $7$. We denote this by $2 \nmid 7$.

The fact that not every integer is a divisor of every other integer is what makes number theory interesting. The question of which integers divide a given integer is of crucial importance as the factorisation of large numbers plays a key role in modern cryptography and secure communications. It can be reduced to the question of finding the prime factorisation of a number. While in general it is hard to do this there is a related problem that is a lot easier to solve, and we will start with that. Let $a, b, d$ be integers.

If $d|a$ and $d|b$ we say that $d$ is a *common divisor* (or *common factor*) of $a$ and $b$; for instance, $1$ is a common divisor of any pair of integers $a$ and $b$.

The *greatest common divisor* (or *highest common factor*) of $a$ and $b$ is the unique integer $d$ satisfying

- $d|a$ and $d|b$ (so that $d$ is a common divisor),
- If $c|a$ and $c|b$ then $c \leq d$ (so that no common divisor exceeds $d$).

We denote the greatest common divisor of $a, b$ by $\gcd(a, b)$.

**Example 3.3.** Find the gcd of $a = 30$ and $b = 42$:

The set of divisors of $30$ is $X = \{\pm 1, \pm 2, \pm 3, \pm 5, \pm 6, \pm 10, \pm 15, \pm 30\}$.

The set of divisors of $42$ is $Y = \{\pm 1, \pm 2, \pm 3, \pm 6, \pm 7, \pm 14, \pm 21, \pm 42\}$.

The set of common divisors of $30$ and $42$ is $X \cap Y = \{\pm 1, \pm 2, \pm 3, \pm 6\}$

This is fine for relatively small integers such as 30 and 42, but how would we find the greatest common divisor of $765432$ and $56789$?

## 3.2   The Euclidean Algorithm

For larger integers we can automate the process using one of the oldest algorithms in mathematics, Euclid's algorithm:

Euclid's algorithm (published in Book VII of *Euclid's Elements* around 300 BC) is based on the following simple observation:

If $a, b$ are integers with $a > b$ then $\gcd(a, b) = \gcd(a - b, b)$.

By repeated application of Euclid's observation, we can reduce the size of the numbers involved in our calculations. For example, suppose we wish to calculate $\gcd(765432, 56789)$. Euclid's observation says that as

$765432 - 56789 = 708643$ then $\gcd(765432, 56789) = \gcd(708643, 56789)$. But then $708643 - 56789 = 651854$ and so $\gcd(765432, 56789) = \gcd(651854, 56789)$. We continue our calculations in this fashion:

$$
\begin{aligned}
765432 - 56789 &= 708643 \\
708643 - 56789 &= 651854 \\
651854 - 56789 &= 595065 \\
&\vdots
\end{aligned}
$$

We can speed this process up slightly by observing that if $a, b$ are integers with $a > b$ then

$$
\begin{aligned}
\gcd(a, b) &= \gcd(a - b, b) \\
&= \gcd(a - 2b, b) \\
&\vdots \\
&= \gcd(r, b)
\end{aligned}
$$

where $r$ is the remainder when we divide $a$ by $b$. Now we can subtract 56789 from 765432, 13 times.

$$
\begin{aligned}
765432 - 13 * (56789) &= 27175 \\
\gcd(765432, 56789) &= \gcd(56789, 27175)
\end{aligned}
$$

We then repeat the process with the two smaller numbers 56789 and 27175.

$$
\begin{aligned}
765432 - (13) * (56789) &= 27175 \\
\gcd(765432, 56789) &= \gcd(56789, 27175) \\
56789 - (2) * (27175) &= 2439 \\
\gcd(765432, 56789) &= \gcd(27175, 2439)
\end{aligned}
$$

Carrying on in this way we can deduce the value of the $\gcd(765432, 56789)$.

$$
\begin{aligned}
765432 - (13) * (56789) &= 27175 \\
56789 - (2) * (27175) &= 2439 \\
27175 - (11) * (2439) &= 346 \\
2439 - (7) * (346) &= 17 \\
346 - (20) * (17) &= 6 \\
17 - (2) * (6) &= 5 \\
6 - (1) * (5) &= 1 \\
5 - (5) * (1) &= 0.
\end{aligned}
$$

So gcd of $765432$ and $56789$ is 1.

# The Remainder Theorem

At each stage of the process above, given integers $a$ and $b$, we have to find integers $q$ and $r$ such that

$$a = qb + r \qquad \text{and} \qquad 0 \le r < b.$$

Recall that Theorem 2.14 which we proved in Chapter 2 tells us that for any **positive** integers $a, b$ there are unique values of $q$ and $r$ satisfying the above conditions.

We will now generalise the result to allow negative values.

In order to be able to apply Lemma 3.12, we require the following important theorem.

---

**Theorem 3.4: Remainder Theorem**

If $a$ and $b$ are integers with $b \ne 0$, then there is a unique pair of integers $q$ and $r$ such that

$$a = qb + r \qquad \text{and} \qquad 0 \le r < |b|.$$

---

Here $|n|$ denotes the absolute value of $n$ which is defined by

$$
|n| = \begin{cases} n & \text{if } n \ge 0 \\ -n & \text{if } n \le 0 \end{cases}
$$

Note that when $n = 0$ we can use either $n$ or $-n$ since $-0 = 0$ (Exercise: Show this from the axioms of $\mathbb{Z}$).

The absolute value is constructed to be non-negative: When $n \geq 0$ then $|n| = n \geq 0$, and when $0 \geq n$ then $|n| = 0 - n \geq n - n = 0$ by *(Sh)*.

*Proof.* As $b \neq 0$ we have $|b| \neq 0$. Since $|b| \geq 0$ for any integer $b$, we have $|b| > 0$.

We know $a \geq 0$ or $a \leq 0$ so we consider cases.

Suppose $a \geq 0$:

Then we can apply Theorem 2.14 to the integers $a \geq 0, |b| > 0$ to show

$$\exists q', r \in \mathbb{Z} \text{ such that } a = q'|b| + r \text{ and } 0 \leq r < |b|.$$

If $b > 0$ then $|b| = b$ and $q = q', r$ satisfy $a = qb + r$ as required.

If $b < 0$ then let $q = -q'$. Then $qb = (-q)(-b) = (-(-q'))|b| = q'|b|$ (see exercises).

So again $a = qb + r$ as required.

Hence the result is true when $a \geq 0$.

Suppose $a \leq 0$:

Then $|a| = -a \geq 0$ so we can apply the first case to show

$$\exists q'', r' \in \mathbb{Z} \text{ such that } -a = q''b + r' \text{ and } 0 \leq r' < |b|.$$

Hence $a = -q''b - r'$.

If $r' = 0$ then $q = -q'', r = 0$ satisfy $a = qb + r$ as required.

If $r' > 0$ then $-r' < 0$ so $|b| - r' < |b|$ while $r' < |b|$ implies $0 < |b| - r'$ by *(Sh)*.

Let $r = |b| - r'$. We have $a = -q''b - |b| + r$.

Letting $q$ be $-q'' - 1$ if $b > 0$ and $-q'' + 1$ if $b < 0$ we have $qb = -q''b - |b|$ so $a = qb + r$ as required. $\qquad\square$

**Exercise 3.5.** Let $a$ and $b$ be integers with $b \neq 0$. Show that there exists a unique pair of integers $q$ and $r$ such that

$$a = qb + r \qquad \text{and } -|b|/2 < r \leq |b|/2.$$

The following example will prove useful in Chapter 5.

**Example 3.6.** If $n$ is a square then $n$ leaves a remainder $0$ or $1$ when divided by $4$.

*Solution.* Let $n = a^2$.

According to Theorem 3.4 $a = 4q + r$ where $r = 0, 1, 2$ or $3$, so that

$$n = (4q + r)^2 = 16q^2 + 8qr + r^2.$$

- If $r = 0$ then $n = 4(4q^2 + 2qr) + 0$,
- if $r = 1$ then $n = 4(4q^2 + 2qr) + 1$,
- if $r = 2$ then $n = 4(4q^2 + 2qr + 1) + 0$,
- and if $r = 3$ then $n = 4(4q^2 + 2qr + 2) + 1$.

This completes the result.

**Example 3.7.** Show that any number that is both a square and a cube leaves a remainder of either 0 or 1 when divided by $7$.

*Solution.* Let $a$ be an integer. By the Remainder Theorem there exists integers $q$ and $r$ such that $a = 7q + r$ and $0 \leq r \leq 6$. Hence

$$a^2 = (7q + r)^2 = 49q^2 + 14qr + r^2 = 7(7q^2 + 2qr) + r^2.$$

Then $r^2 \in \{0^2, \ldots, 6^2\} = \{0, 1, 4, 9, 16, 25, 36\}$. But $9 = 7 + 2$, $16 = 7 \times 2 + 2$, $25 = 7 \times 3 + 4$, $36 = 7 \times 5 + 1$ and so $a^2 = 7q' + r'$ where $r' \in \{0, 1, 2, 4\}$.

Using a similar argument, we see that if $b = 7q + r$ then

$$b^3 = (7q + r)^3 = 343q^3 + 147q^2r + 21qr^2 + r^3.$$

We can then easily check that $r^3$ leaves a remainder of $0, 1$ or $6$ on division by 7 and hence $b^3 = 7q'' + r''$ where $r'' \in \{0, 1, 6\}$.

So if a number is both a square and a cube it must have a remainder of either 0 or 1 when divided by $7$.

We shall consider lots of examples of Euclid's algorithm later in this chapter. For now, we want to focus on proving that Euclid's observation is true and that the resultant algorithm gives the correct result.

## Recall our definition of a divisor

We say that an integer $d$ is a divisor of the integer $a$ if and only if there is an integer $b$ such that $a = db$.

Notice that the definition is phrased entirely using properties of the integers encoded in the axioms given in Chapter 2.

Now what can we say about divisors using the axioms and the definitions?

- Every integer divides $0$. This is because $0 = d \cdot 0$. Notice that this is not one of our axioms but we can deduce it from the axioms, and will do so in one of the tutorial sheets for the course.
- $1$ divides every integer. This is because $a = 1 \cdot a$ from axiom *(Id)*.
- Every integer divides itself. This is because $a = a \cdot 1$ from axiom *(Id)*.
- $0$ does not divide any integer except itself. If $a = 0 \cdot b$ then $a = 0$, from the first of these observations.

We have already mentioned that we shall attempt to avoid using division whenever possible. The following important result will often help us in this task.

---

**Lemma 3.8: Cancellation Lemma**

If $m, n, p$ are integers with $m \cdot p = n \cdot p$ then either $p = 0$ or $m = n$.

---

*Proof.* Since $m \cdot p = n \cdot p$ we have $m \cdot p - n \cdot p = n \cdot p - n \cdot p = 0$ by axiom *(Neg)*.

By axiom *(Comm)* we get $p \cdot m - p \cdot n = 0$ and by axiom *(Dist)* we get $p(m-n) = 0$.

Now by axiom *(ZD)* either $p = 0$ or $m - n = 0$, and so $p = 0$ or $(m-n)+n = 0+n$ which means $m + (-n + n) = n$ by *(Id)* and *(As)*. Hence $m = m + 0 = m + (n - n) = n$ by *(Id)*, *(Neg)* and *(Comm)* as required. □

---

**Lemma 3.9**

Let $a, b$ and $c$ be integers.

1. If $a|b$ and $b|c$ then $a|c$.
2. If $a|b$ and $c|d$ then $ac|bd$.
3. If $m \neq 0$, then $a|b$ if and only if $ma|mb$.
4. If $d|a$ and $a \neq 0$ then $|d| \leq |a|$.

---

*Proof.* Let $a, b$ and $c$ be integers.

1. If $b = ma$ for some integer $m$, and $c = nb$ for some integer $n$, then $c = n(ma) = (nm)a$ and $nm$ is an integer.

2. If $b = ma$ for some integer $m$ and $d = nc$ for some integer $n$, then $bd = (ma)(nc) = ((ma)n)c = (m(an))c = (m(na))c = ((mn)a)c = (mn)(ac)$.

3. If $ma|mb$ then there is an integer $k$ such that $mb = k(ma) = m(ka)$. Since $m \neq 0$ we can apply the cancellation lemma to see that $ka = b$. It follows that $a|b$ as required.

   Now suppose that $a|b$. Since $m|m$ and $a|b$ it follows that $ma|mb$ by part 2.

4. If $a = md$ for some integer $m$ then

$$|a| = |md| = |m||d|.$$

   Since $a \neq 0$ then $m \neq 0$ and so $|m| \geq 1$. Hence $|a| \geq |d|$.

This concludes the proof.                                                    □

Notice that in part (4), we have mentioned that $|md| = |m||d|$. This is not one of the axioms, but can be deduce from the axioms and will appear as an example in one of the tutorial sheets for the course. Let $a_1, \ldots, a_k, u_1, \ldots, u_k$ be integers. We shall refer to an expression of the form

$$a_1 u_1 + \cdots + a_k u_k$$

as a *linear combination* of the integers $a_1, \ldots, a_k$.

---

**Theorem 3.10**

Let $a_1, \ldots, a_k, u_1, \ldots, u_k, a, b, c$ be integers

1. If $c$ divides $a_1, \ldots, a_k$, then $c$ divides any linear combination of the integers $a_1, \ldots, a_k$.
2. $a|b$ and $b|a$ if and only if $a = \pm b$.

---

*Proof.* Let $a_1, \ldots a_k, u_1, \ldots u_k, a, b, c$ be integers

1. If $c$ divides $a_i$ then $a_i = q_i c$ for some integers $q_i$ ($i = 1, \ldots, k$). Then $a_1 u_1 + \cdots + a_k u_k = q_1 c u_1 + \cdots + q_k c u_k = (q_1 u_1 + \cdots + q_k u_k)c$, and as $q_1 u_1 + \cdots + q_k u_k \in \mathbb{Z}$ (since $q_i, u_i \in \mathbb{Z}$) we see that $c|(a_1 u_1 + \cdots + a_k u_k)$.

NOTE:

The fact that we can "pull" the factor of $c$ out of this expression is due to the distributive law *(Dist)*, the commutative law *(Comm)* and the associative law *(As)* for integer multiplication and addition.

2. If $a = \pm b$ then $b = qa$ and $a = q'b$ where $q = q' = \pm 1$. Hence $a|b$ and $b|a$. Conversely, suppose that $a|b$ and $b|a$, so $b = qa$ and $a = q'b$ for some

integers $q$ and $q'$.

If $b = 0$ then the second equation gives $a = 0$, and so $a = \pm b$ as required. Suppose now that $b \neq 0$.

Combining the two equations, we have $1.b = b = qq'b$, and so by the cancellation lemma, and since $b \neq 0$, we deduce that $qq' = 1$.

Hence $q$ and $q'$ both divide $1$, and by Lemma 3.9(4), $|q| \leq |1| = 1$ and $|q'| \leq 1$.

Hence $q$ and $q'$ are both in the set $\{-1, +1\}$ and $a = \pm b$.

$\square$

We get the following very useful Corollary for the case $k = 2$.

> **Corollary 3.11**
>
> If $c$ divides $a$ and $b$, then $c$ divides $au + bv$ for all integers $u$ and $v$.

We are now in a position to prove Euclid's observation mentioned earlier.

> **Lemma 3.12**
>
> If $a, b$ are integers with $a = qb + r$ then $\gcd(a, b) = \gcd(b, r)$.

*Proof.* By Corollary 3.11, any common divisor of $b$ and $r$ also divides any integer of the form $ub + vr$.

Now if $a = qb + r$ we can put $u = q$ and $v = 1$ to see that *any* common divisor of $b$ and $r$ is also a divisor of $a$. Hence any common divisor of $b$ and $r$ is also a common divisor of $a$ and $b$.

Conversely, any common divisor of $a$ and $b$ also divides $a - qb$ (again by Corollary 3.11, with $u = 1$ and $v = -q$).

Hence any common divisor of $a$ and $b$ is also a common divisor of $b$ and $r = a - qb$.

Thus the two pairs $a, b$ and $b, r$ have the same common divisors, so they have the same greatest common divisor. $\square$

## Proof of Euclid's Algorithm

We are now in a position to prove that Euclid's algorithm, both terminates, and produces the desired greatest common divisor. We use the remainder theorem to successively define the integers $r_i$ so that $r_1$ is the remainder upon dividing $a$ by $b$ and for $i \geq 1$ each $r_{i+1}$ is the remainder upon dividing $r_{i-1}$ by $r_i$.

$$\begin{array}{rcll}
a & = & q_1 b + r_1 & 0 \leq r_1 < b \\
b & = & q_2 r_1 + r_2 & 0 \leq r_2 < r_1 \\
r_1 & = & q_3 r_2 + r_3 & 0 \leq r_3 < r_2 \\
& \cdots & & \\
r_{i-1} & = & q_{i+1} r_i + r_{i+1} & 0 \leq r_{i+1} < r_i \\
& \cdots & & \\
r_{n-1} & = & q_{n+1} r_n + r_{n+1} & 0 \leq r_{n+1} < r_n \\
r_n & = & q_{n+2} r_{n+1} + 0. &
\end{array}$$

In particular we have $a \geq b > r_1 > r_2 > \ldots \geq 0$. By the well ordering principle this process stops, and by inspection we see it can only stop when the final remainder is $0$.

By Lemma 3.12 we see that

$$\gcd(a, b) = \gcd(b, r_1) = \gcd(r_1, r_2) = \ldots$$

and in general each pair $r_i, r_{i+1}$ has the same greatest common divisor as the successor pair $r_{i+1}, r_{i+2}$ so in particular (and by induction) $\gcd(a, b) = \gcd(r_n, r_{n+1}) = r_{n+1}$ as required.

**Example 3.13.** To calculate $d = \gcd(1815, 1415)$ we write

$$\begin{array}{rcl}
1815 & = & 1 \times 1415 + 400 \\
1415 & = & 3 \times 400 + 215 \\
400 & = & 1 \times 215 + 185 \\
215 & = & 1 \times 185 + 30 \\
185 & = & 6 \times 30 + 5 \\
30 & = & 6 \times 5 + 0
\end{array}$$

The last non-zero remainder is $5$, so $d = 5$.

In many cases, the value of $d$ can be identified before a zero remainder is reached: since $d = \gcd(a, b) = \gcd(b, r_1) = \gcd(r_1, r_2) = \ldots$, one can stop as soon as one recognises the greatest common divisor of a pair of consecutive terms in the sequence $a, b, r_1, r_2, \ldots$. In Example 3.13, for instance, the remainders $185$ and $30$ clearly have greatest common divisor $5$, so $d = 5$.

**Exercise 3.14.** Use Exercise 3.5 to devise an alternative algorithm to Euclid's for calculating the greatest common divisor. This is known as the *least remainders algorithm*.

# Euclid's algorithm for polynomials

A polynomial over the integers is a function like $5x^3 + 2x^2 + 7x - 4$ with integer coefficients. It's degree is the highest power of the variable which arises in the expression. [There is a much more formal definition which can be given, but it agrees with this naive idea.]

The set of all integer polynomials shares certain similarities with the integers. In particular we can add, subtract and multiply polynomials just as we can integers. Note that the polynomials include **constant** polynomials, i.e. the integers.

The operations on polynomials satisfy all of the algebra axioms – *(Op),(Id),(Neg),(Comm),(As),(Dist)* and *(ZD)* – as the integers do.

To have a version of the Remainder Theorem we also need the idea of ordering. One approach is to say that $f \leq g$ if and only if $f = g$ or $f$ has strictly lower degree than $f$. The idea here is that "smaller" polynomials are ones of lower degree. This order satisfies *(Ref),(AS),(Tr),(Sc)* and a variation of *(WO)*[2].

The idea is then to apply a version of the remainder theorem to polynomials $f, g$ (essentially long division of polynomials) to write

$$g = q \cdot f + r$$

where $q, r$ are polynomials and $r < f$ (meaning that the *degree* of $r$ is less than the degree of $f$).

Lemma 3.12 works as before to say that the common divisors of $q, f$ are the same as the common divisors of $f, r$, to these have the same greatest common divisors (though now there may be more than one "greatest" common divisor since greatest is defined only in terms of the degree of the polynomials).

The question of divisibility of integers complicates the question of divisibility of integer polynomials. A version of the remainder theorem for polynomials will work provided that the leading coefficient[3] of $f$ is a divisor of the leading coefficient of $g$. Things are much simpler if we look at polynomials with rational coefficients in which case the remainder theorem always works.

---

[2]Non-empty sets of polynomials will have *one or more* least elements.
[3]The leading coefficient of a polynomial is the coefficient of the highest power of $x$.

You will find a description of Euclid's algorithm in the context of polynomials on the page

http://en.wikipedia.org/wiki/Greatest common divisor of two polynomials

## Rings

**Definition 3.15**

A *ring* is a set $R$ equipped with two binary operations $+, \cdot$ and satisfying the axioms *(Op),(Id),(Neg),(Comm),(As),(Dist)*.

- Of course $\mathbb{Z}$ is an example of a ring as we based the definition of ring on (most of) the algebraic axioms of $\mathbb{Z}$.
- The set of polynomials with coefficients in $\mathbb{Z}$ is a ring.
- The rational numbers $\mathbb{Q}$, the real numbers $\mathbb{R}$ and the complex numbers $\mathbb{C}$ are examples of rings.
- The set of polynomials with coefficients in $\mathbb{Q}, \mathbb{R}$ or $\mathbb{C}$ is a ring.
- The set $\{0, 1\}$ with the operations $\oplus$ defined by $0 \oplus 0 = 1 \oplus 1 = 0$ and $0 \oplus 1 = 1 \oplus 0 = 1$ and the usual multiplication $\cdot$ is a ring. We can think of this example as representing even/odd (where $0$ is interpreted as meaning even and $1$ meaning odd):

  $$\text{even} \oplus \text{even} = \text{odd} \oplus \text{odd} = \text{even}, \qquad \text{even} \oplus \text{odd} = \text{odd} \oplus \text{even} = \text{odd}$$

We will see other examples of rings in later chapters when we study *modular arithmetic*.

## 3.3　Bezout's Identity

**Theorem 3.16: Bezout's Identity**

If $a$ and $b$ are integers (not both $0$), then there exist integers $u$ and $v$ such that
$$\gcd(a, b) = au + bv.$$
Moreover $\gcd(a, b)$ is the least positive integer of the form $au + bv$ ($u, v \in \mathbb{Z}$).

NOTE: *The values of $u$ and $v$ are not uniquely determined by $a$ and $b$.*

The proof of Bezout's Identity will involve sets called *ideals*.

---
**Definition 3.17**

Let $R$ be a *ring*. A subset $I$ of $R$ is an *ideal* if:

- $I$ is non-empty;
- for all $a, b$ in $I$, $a + b$ is also in $I$
- for all $a \in I$ and $r \in R$, $ra$ is also in $I$

---

Note the last requirement says $I$ is closed under multiplication by **all** elements of $R$ not just under multiplication by other elements of $I$.

We consider the ring $\mathbb{Z}$:

- The set $I = \{0\}$ is an ideal of $\mathbb{Z}$.
- The set of *even numbers* $I = \{\ldots, -2, 0, 2, 4, \ldots\}$ is an ideal of $\mathbb{Z}$.
- The set of *odd numbers* is not an ideal (it is not closed under addition, nor is it closed under multiplication by integers).
- The set $\mathbb{N}_0$ of non-negative integers is not an ideal (it *is* closed under addition, but not under multiplication by integers).

---
**Lemma 3.18**

Let $a, b \in \mathbb{Z}$. Then the set

$$I = \{au + bv : u, v \in \mathbb{Z}\}$$

is an ideal in $\mathbb{Z}$.

---

Note: The set $I$ contains both $a$ and $b$: setting $u = 1$ and $v = 0$ we have $au + bv = a$, while setting $u = 0, v = 1$ we get $au + bv = b$.

*Proof.* To show $I$ is an ideal we must prove three things.

$I$ is non-empty:

Let $u = v = 0$. Then $au + bv = 0$ (by Lemma 2.3) so $0 \in I$.

In particular $I$ is non-empty.

For $x, y$ in $I$, $x + y$ is also in $I$:

If $x \in I$ then there exists $u, v$ in $\mathbb{Z}$ such that $x = au + bv$.

If $y \in I$ then there exists $u', v'$ in $\mathbb{Z}$ such that $y = au' + bv'$.

So $x + y = au + bv + au' + bv' = (au + au') + (bv + bv')$ by *(Comm),(As)*.

Hence $x + y = a(u + u') + b(v + v')$ *(Dist)*.

So $x + y = au'' + bv''$ where $u'' = u + u'$ and $v'' = v + v'$ are in $\mathbb{Z}$.

Therefore $x + y \in I$.

For $x$ in $I$, $r \in \mathbb{Z}$, $rx$ is also in $I$:

If $x \in I$ then there exists $u, v$ in $\mathbb{Z}$ such that $x = au + bv$.

So $rx = r(au + bv) = rau + rbv$ by *(Dist)*.

This is $a(ru) + b(rv)$ by *(As),(Comm)* so $rx = au' + bv'$ where $u' = ru, v' = rv$ are integers.

Hence $rx \in I$.                                                                    □

---

**Lemma 3.19**

Let $I$ be a non-zero ideal in $\mathbb{Z}$. Then there exists $d \in I$ with $d$ positive such that
$$x \in I \iff d|x$$
and $d$ is the least positive element of $I$.

---

Note: this means that $I$ has the form:
$$I = \{nd : n \in \mathbb{Z}\}$$

*Proof.* As $I$ is non-zero it contains some $x \neq 0$. As $I$ is an ideal it must also contain $(-1)x = -x$.

Hence $I$ contains some positive element so the set $\{x \in I : x > 0\}$ has a least element by *(WO)*.

Let $d$ be the least element of this set, i.e. $d$ is the least positive element of $I$.

We will prove: 1. If $x = nd$ for some $n \in \mathbb{Z}$ then $x \in I$. 2. If $x \in I$ then f $x = nd$ for some $n \in \mathbb{Z}$.

1. As $d \in I$, if $x = nd$ with $n \in \mathbb{Z}$ then $x \in I$ by the definition of ideal.

2. If $x \in I$ then, by the remainder theorem 3.4 there exists $q, r \in \mathbb{Z}$ with $x = qd + r$ and $0 \leq r < d$.

Since $d \in I$ and $-q \in \mathbb{Z}$ we have $-(qd) = (-q)d \in I$ by the definition of ideal.

By assumption $x \in I$ and so as $qd \in I$ we have $r = x - qd$ in $I$.

By definition, $d$ is the **least** positive element of $I$ and $r < d$ hence $r$ cannot be positive.

But $r \geq 0$ so $r = 0$. Thus $x = qd$ as required. □

---

## Proof of Bezout's Identity

*Proof.* Let $I = \{au + bv : u, v \in \mathbb{Z}\}$. This is an ideal by Lemma 3.18.

So by Lemma 3.19 there is a positive $d \in I$ such that $d|x$ for all $x \in I$.

As $d \in I$ there exist integers $u, v$ such that $d = au + bv$.

As $a, b \in I$ we have $d|a$ and $d|b$, that is $d$ is a common factor of $a, b$.

If $c$ is any other common divisor of $a, b$ then $c$ divides $au + bv = d$ by Corollary 3.11.

Hence $|c| \leq |d|$ by Lemma 3.9. Since $|d| = d$ we have $c \leq |c| \leq d$ for any common divisor $c$ of $a, b$.

Thus $d$ is the greatest common divisor of $a, b$ so $\gcd(a, b) = d = au + bv$. □

The theorem tells us that $\gcd(a, b) = au + bv$ has integer solutions, but, like a person in a desert looking for an oasis, knowing that it exists is not as useful as knowing how to find it!

We can use the equations which arise when we apply Euclid's algorithm to calculate $d = \gcd(a, b)$:

The penultimate equation, of Euclid's algorithm has the form

$$r_{n-1} - q_{n+1}r_n = r_{n+1},$$

Recall that $d = r_{n+1}$. So we get $d = r_{n-1} - q_{n+1}r_n$. Marching back up Euclid's algorithm we have:

$$r_n = r_{n-2} - q_nr_{n-1},$$

This allows us to eliminate $r_n$

$$d = r_{n-1} - q_{n+1}(r_{n-2} - q_nr_{n-1}) = (1 + q_{n-1}q_n)r_{n-1} - q_{n+1}r_{n-2},$$

- Gradually work backwards through the equations in the algorithm, eliminating $r_{n-1}, r_{n-2}, \ldots$ in succession.
- eventually we express $d$ in terms of $r_1$ and $r_2$,
- then in terms of $b$ and $r_1$,
- and finally in terms of $a$ and $b$.

**Example 3.20.** In Example 3.13 we used Euclid's algorithm to calculate $d$, where $a = 1815$ and $b = 1415$. Use those equations again to obtain $u$ and $v$.

$$
\begin{aligned}
5 &= 1 \times 185 - 6 \times 30 \\
&= 1 \times 185 - 6 \times (215 - 1 \times 185) \\
&= 7 \times 185 - 6 \times 215 \\
&= 7 \times (400 - 1 \times 215) - 6 \times 215 \\
&= 7 \times 400 - 13 \times 215 \\
&= 7 \times 400 - 13 \times (1414 - 3 \times 400) \\
&= 46 \times 400 - 13 \times 1414 \\
&= 46 \times (1815 - 1 \times 1415) - 13 \times 1415 \\
&= 46 \times 1815 - 59 \times 1415
\end{aligned}
$$

and so $u = 46, v = -59$.

**Example 3.21.** Find $d = \gcd(1068, 4294)$ and express $d$ in the form $d = 1068u + 4294v$:

First we use Euclid's algorithm

$$
\begin{aligned}
4294 &= 4 \times 1068 + 22 \\
1068 &= 48 \times 22 + 12 \\
22 &= 1 \times 12 + 10 \\
12 &= 1 \times 10 + 2 \\
10 &= 5 \times 2 + 0
\end{aligned}
$$

Next we apply Bezout's identity

$$
\begin{aligned}
2 &= 12 - 10 \\
2 &= 12 - (22 - 12) \\
2 &= 2 \times 12 - 22 \\
2 &= 2 \times (1068 - 48 \times 22) - 22 \\
2 &= 2 \times 1068 - 97 \times 22 \\
2 &= 2 \times 1068 - 97(4294 - 4 \times 1068) \\
2 &= 390 \times 1068 - 97 \times 4294
\end{aligned}
$$

> **Theorem 3.22: Lamé's Theorem**
>
> The number of steps in the Euclidean algorithm does not exceed 5 times the number of decimal digits in the smaller of the two numbers.

*Proof.* Writing $a = r_0$ and $b = r_1$ we can apply the Euclidean algorithm to $a$ and $b$ to get a sequence of equations.

$$
\begin{aligned}
r_0 &= q_1 r_1 + r_2 & 0 \leq r_2 < r_1 \\
r_1 &= q_2 r_2 + r_3 & 0 \leq r_3 < r_2 \\
r_2 &= q_3 r_3 + r_4 & 0 \leq r_4 < r_3 \\
&\phantom{=}\cdots \\
r_{n-3} &= q_{n-2} r_{n-2} + r_{n-1} & 0 \leq r_{n-1} < r_{n-2} \\
r_{n-2} &= q_{n-1} r_{n-1} + r_n & 0 \leq r_n < r_{n-1} \\
r_{n-1} &= q_n r_n
\end{aligned}
$$

Notice that $q_1, \ldots, q_{n-1} \geq 1$ and that $q_n \geq 2$ (as $r_n < r_{n-1}$). Consequently we see that

$$
\begin{aligned}
r_n \geq 1 &= f_2 \\
r_{n-1} \geq 2 r_n \geq 2 f_2 &= f_3 \\
r_{n-2} \geq r_{n-1} + r_n \geq f_3 + f_2 &= f_4 \\
&\cdots \\
r_2 \geq r_3 + r_4 \geq f_{n-1} + f_{n-2} &= f_n \\
b = r_1 \geq r_2 + r_3 \geq f_n + f_{n-1} &= f_{n+1}
\end{aligned}
$$

But from Example 2.13 we see that $f_{n+1} \geq \alpha^{n-1}$ for $n \geq 2$, where $\alpha = (1 + \sqrt{5})/2$. Now since $\log_{10} \alpha > 1/5$ it follows that

$$
\log_{10} b \geq (n-1) \log_{10} \alpha > (n-1)/5.
$$

If $b$ has $k$ decimal digits then $b < 10^k$ and so $\log_{10} b < k$. Hence $n - 1 < 5k$ and the result follows. $\qquad\square$

So, for example, if we apply the Euclidean algorithm to integers with 100 digits, then we require no more than 500 steps in the algorithm.

## The gcd of a set of integers

Having seen how to calculate the greatest common divisor of two integers, it is a straightforward matter to extend this to any finite set of integers (not all $0$). The method, which involves repeated use of Euclid's algorithm, is based on the following observation.

**Exercise 3.23.** Prove that for any $k \geq 2$,

$$
\gcd(a_1, \ldots, a_k) = \gcd(\ldots \gcd(\gcd(a_1, a_2), a_3), \ldots, a_k).
$$

For example, to calculate $\gcd(36, 24, 54, 27)$ we calculate

$$
\begin{aligned}
d_2 &= \gcd(36, 24) &=& 12 \\
d_3 &= \gcd(12, 54) &=& 6 \\
\text{and finally } d_4 &= \gcd(6, 27) &=& 3
\end{aligned}
$$

**Example 3.24.** We saw in Example 3.13 that if $a = 1815$ and $b = 1415$ then $d = 5$, so the integers of the form $c = 1815x + 1415y$ are the multiples of $5$. Example 3.20 gives $5 = 1815 \times 46 + 1415 \times (-59)$, so multiplying through by $e$ we can express any multiple of 5, $5e$, in the form $1815x + 1415y$: for instance, $-10 = 1815 \times (-92) + 1415 \times 118$.

## 3.4   Coprime Integers

Two integers $a$ and $b$ are *coprime* (or *relatively prime*) if $\gcd(a, b) = 1$.

For example, $15$ and $23$ are coprime, but $15$ and $20$ are not.

When dealing with more than two integers, there is two ways to generalise.

- The integers $a_1, a_2, \ldots, a_n$ are *coprime* if $\{\gcd(a_1, a_2, \ldots, a_n) = 1.\}$
- The integers $a_1, a_2, \ldots, a_n$ are *mutually coprime* if $\{\gcd(a_i, a_j) = 1$ whenever $i \neq j.\}$

It should be clear that if the integers $a_1, a_2, \ldots, a_n$ are mutually coprime then they are coprime (since $\gcd(a_1, a_2, \ldots) | \gcd(a_i, a_j)$).

*The converse is false: for example, the integers* $6, 10, 15$ *are coprime but are not mutually coprime.*

> **Corollary 3.25**
>
> Two integers $a$ and $b$ are coprime if and only if there exist integers $x$ and $y$ such that
> $$ax + by = 1.$$

*Proof.* If $\gcd(a, b) = 1$ then by Theorem 3.16, we have $ax + by = 1$ for some $x, y \in \mathbb{Z}$.

Conversely we know that $\gcd(a, b)$ is the least positive value of $ax + by$, so if this is 1 for some $x, y$ then $\gcd(a, b) = 1$. $\qquad\square$

The following two corollaries will prove extremely useful, and will be used on a number of occasions in subsequent chapters.

> **Corollary 3.26**
>
> If $\gcd(a, b) = d$ then
> $$\gcd(ma, mb) = md$$
> for every integer $m > 0$, and
> $$\gcd(u, v) = 1$$
> where $u$ is the unique integer such that $a = ud$ and $v$ is the unique integer such that $b = vd$.

*Proof.* By Theorem 3.16, $\gcd(ma, mb)$ is the smallest positive value of $max + mby = m(ax + by)$, where $x, y \in \mathbb{Z}$, while $d$ is the smallest positive value of $ax + by$, so $\gcd(ma, mb) = md$. Write $d = ax + by = (du)x + (dv)y$. Then by the cancellation lemma $ux + vy = 1$. Corollary 3.25 then implies that $u$ and $v$ are coprime. $\qquad\square$

> **Corollary 3.27**
>
> Let $a$ and $b$ be coprime integers.
>
>   a. If $a|c$ and $b|c$ then $ab|c$.
>   b. If $a|bc$ then $a|c$.

*Proof.* Let $a$ and $b$ be coprime integers.

   a. We have $ax + by = 1$, $c = ae$ and $c = bf$ for some integers $x, y, e$ and $f$. Then $c = cax + cby = (bf)ax + (ae)by = ab(fx + ey)$, so $ab|c$.

   b. As in (a), $c = cax + cby$. Since $a|bc$ and $a|a$, Corollary 3.11 implies that $a|(cax + cby) = c$.

$\qquad\square$

## Least Common Multiple

If $a$ and $b$ are integers, then a *common multiple* of $a$ and $b$ is an integer $c$ such that $a|c$ and $b|c$. If $a$ and $b$ are both non-zero, then they have positive common multiples (such as $|ab|$), so by the well-ordering principle they have a *least common multiple* or, more precisely, a least *positive* common multiple; this is the unique positive integer $l$ satisfying

   1. $a|l$ and $b|l$ (so $l$ is a common multiple), and

2. if $a|c$ and $b|c$, with $c > 0$, then $l \le c$ (so no positive common multiple is less than $l$).

> **Theorem 3.28**
>
> Let $a$ and $b$ be positive integers, with $d = \gcd(a, b)$ and $l = \operatorname{lcm}(a, b)$. Then
> $$dl = ab.$$

NOTE: We can assume that $a, b > 0$ since $\gcd(a, b) = \gcd(|a|, |b|)$ and $\operatorname{lcm}(a, b) = \operatorname{lcm}(|a|, |b|)$.

*Proof.* Let $a = ed$ and $b = fd$, so that
$$ab = efd^2.$$

Now consider $efd$, clearly this is positive, so we can show that it is equal to $l$ by showing that it satisfies conditions (1) and (2) of the definition of $\operatorname{lcm}(a, b)$.

Firstly,
$$efd = (ed)f = af \quad \text{and} \quad edf = (fd)e = be;$$
thus $a|def$ and $b|def$, so (1) is satisfied.

Secondly, suppose that $a|c$ and $b|c$, with $c > 0$; it follows that there exist integers $p, q$ such that $c = ap = bq$. We need to show that $def \le c$.

By Theorem 3.16 there exist integers $u$ and $v$ such that $d = au + bv$. Now
$$(ab)(qu + pv) = abqu + abpv = acu + bcv = c(au + bv) = cd$$

hence $ab$ is a divisor of $cd$.

This implies that $cd = kab$ for some integer $k$, so $cd = kd(def)$. The cancellation lemma means that $c = kdef$. Hence $def$ divides $c$, and so, as we noted in Lemma 3.9, $def \le c$ as required. □

**Example 3.29.** If $a = 16$ and $b = 12$, then $d = 4$ and $l = 48$; thus $dl = 192 = ab$, agreeing with Theorem 3.28.

**Example 3.30.** We can use Theorem 3.28 to find $l = \operatorname{lcm}(a, b)$ efficiently by first using Euclid's algorithm to find $d = \gcd(a, b)$, and then calculating $l = ab/d$.

Since $\gcd(1815, 1415) = 5$ we have $\operatorname{lcm}(1815, 1415) = (1815 \times 1415)/5 = 513645$.

## 3.5  Linear Diophantine Equations

Find the general solution to

$$1815x + 1415y = -10,$$

Equations in one or more variables, for which we seek integer-valued solutions.

Simplest example is the *linear Diophantine equation*

$$ax + by = c$$

The following result was known to the Indian mathematician Brahmagupta, around 628 AD:

---

**Theorem 3.31**

Let $a, b$ and $c$ be integers, with $a$ and $b$ not both $0$, and let $d = \gcd(a, b)$. Then the equation

$$ax + by = c$$

has an integer solution $x, y$ if and only if $c$ is a multiple of $d$, in which case there are infinitely many solutions.

Writing $a = pd$ and $b = qd$ the solutions are

$$x = x_0 + qn, \quad y = y_0 - pn \qquad (n \in \mathbb{Z}),$$

where $x_0, y_0$ is any particular solution.

---

*Proof.* The fact that there is a solution if and only if $d|c$ is simply a restatement of Theorem 3.16.

To find the general form of the solutions, let $x_0, y_0$ be a particular solution, so $ax_0 + by_0 = c$.

Let $n \in \mathbb{Z}$ and put

$$x = x_0 + qn, \quad y = y_0 - pn.$$

Then

$$ax + by = a(x_0 + qn) + b(y_0 - pn) =$$

$$ax_0 + by_0 + (aqn - bpn) = c + (pdqn - qdpn) = c,$$

and so $x, y$ is also a solution. As this holds for an integer $n$ then there are therefore infinitely many solutions.

Now suppose that $x, y$ is any integer solution, so that $ax + by = c$. Since $ax + by = c = ax_0 + by_0$ we have

$$a(x - x_0) = b(y_0 - y).$$

Now $a = pd$ and $b = qd$ and so

$$pd(x - x_0) = qd(y_0 - y).$$

Applying the cancellation lemma to cancel $d$ we get

$$p(x - x_0) = -q(y - y_0). \qquad (*)$$

Hence $q$ is a divisor of $p(x - x_0)$, and $p$ is a divisor of $q(y - y_0)$.

Since $a$ and $b$ are not both 0, we can suppose that $b \neq 0$ (otherwise we can interchange the roles of $a$ and $b$). Consequently $q \neq 0$.

We also note that $q$ is coprime to $p$ by Corollary 3.26, and so it divides $x - x_0$ by Corollary 3.27(b). Thus $x - x_0 = qn$ for some integer $n$ and so

$$x = x_0 + qn.$$

Substituting back for $x - x_0$ in $(*)$ we get

$$-q(y - y_0) = p(x - x_0) = pqn,$$

Since $q \neq 0$ the Cancellation Lemma implies that

$$y = y_0 - pn.$$

$\square$

Thus we can find the solutions of any linear Diophantine equation $ax + by = c$ by the following technique:

- Calculate $d = \gcd(a, b)$, either directly or by Euclid's algorithm.
- Check whether $d$ is a divisor of $c$: if it is not, there are no solutions, so stop here; if it is, write $c = de$.
- If $d|c$, use the method of proof of Theorem 3.16 to find integers $u$ and $v$ such that $au + bv = d$; then $x_0 = ue, y_0 = ve$ is a particular solution of $ax + by = c$.
- Use Theorem 3.31 to find the general solution $x, y$ of the equation.

**Example 3.32.** Find the general solution to

$$1815x + 1415y = -10,$$

so $a = 1815$, $b = 1415$ and $c = -10$.

In step (1), we use Example 3.13 to see that $d = 5$.

In step (2) we check that $d$ divides $c$: in fact, $c = -2d$, so $e = -2$.

In step (3) we use Example 3.20 to write $d = 46 \times 1815 + (-59) \times 1415$; thus $u = 46$ and $v = -59$,

so $x_0 = 46 \times (-2) = -92$ and $y_0 = (-59) \times (-2) = 118$ give a particular solution of the equation.

By Theorem 3.31, the general solution has the form

$$x = -92 + \frac{1415n}{5} = -92 + 283n, \quad y = 118 - \frac{1815n}{5} = 118 - 363n$$

$$n \in \mathbb{Z}.$$

# Chapter 4

# Prime Numbers

In this chapter we will examine:

- The Fundamental Theorem of Arithmetic
- Euclid's proof that there are infinitely many prime numbers
- The Prime Number Theorem.

The first two results will be proved, but the last, concerning the distribution of primes will only be discussed.

- An integer $p > 1$ is said to be *irreducible* if the only positive divisors of $p$ are $1$ and $p$ itself.
- It is said to be *prime* if for all integers $a$ and $b$, whenever $p|ab$ either $p|a$ or $p|b$.

Note that (by definition) $1$ is not irreducible. The smallest irreducible is $2$, and all the other irreducibles (such as $3, 5, 7, 11, \ldots$) are odd.

An integer $n > 1$ which is not irreducible (such as $4, 6, 8, 9, \ldots$) is said to be *composite*; such an integer has the form $n = ab$ where $1 < a < n$ and $1 < b < n$.

In school the definition of a prime that is given actually coincides with our definition of an irreducible number. Our first result shows that this confusion is justified:

> **Lemma 4.1**
>
> Let $a, p$ be integers.
>
>   1. if $p$ is irreducible, either $p$ divides $a$, or $a$ and $p$ are coprime;
>   2. if $p$ is irreducible then it is also prime, i.e., if $p$ divides $ab$, then $p$ divides $a$ or $p$ divides $b$.
>   3. if $p$ is prime then it is irreducible.

*Proof.* Let $a, p$ be integers.

  1. As $p$ is irreducible and since $\gcd(a, p)$ is a positive divisor of $p$, then $\gcd(a, p)$ must be 1 or $p$. If $\gcd(a, p) = p$, then since it follows that $p|a$ since $\gcd(a, p)$ divides $a$. On the other hand, if $\gcd(a, p) = 1$ then $a$ and $p$ are coprime.
  2. Suppose $p|ab$ and that $p$ does not divide $a$. Then by part (1) $\gcd(a, p) = 1$. Hence by Corollary 3.27(b) it follows that $p|b$ as required.
  3. Suppose that $p$ is prime and that $p = ab$. Notice that we can assume that $a, b > 0$. Since $p$ is prime then $p|a$ or $p|b$. Suppose, without loss of generality, that $p|a$ so that $a = kp$ for some integer $k$. Then $p = kpb$ and hence by the cancellation lemma, it follows that $kb = 1$ and so $b \leq 1$. But $b > 0$ and so $b = 1$, and $p$ is irreducible as required.

$\square$

> **Corollary 4.2**
>
> If $p$ is prime and $p$ divides $a_1 \ldots a_k$, then $p$ divides $a_i$ for some $i$.

*Proof.* We will carry out (weak) induction on $k$. If $k = 1$ then the assumption is that $p|a_1$, so the conclusion is automatically true (with $i = 1$).

Now assume, by way of induction, that $k \geq 1$ and that the result is proved for all products of $k$ factors $a_i$. Let $a = a_1 \ldots a_k$ and $b = a_{k+1}$ so that $a_1 \ldots a_{k+1} = ab$ and then $p|ab$. Since $p$ is prime it follows that $p|a$ or $p|b$. In the first case we have $p|a_1 \ldots a_k$ and the induction hypothesis implies that $p|a_i$ for some $i = 1, \ldots, k$. Otherwise we have $p|a_{k+1}$. Thus in either case $p|a_i$ for some $i$, as required. $\square$

We now come to one of the most important theorems in elementary number theory.

> ### Theorem 4.3: The Fundamental Theorem of Arithmetic
>
> Each integer $n > 1$ has a prime-power factorisation
>
> $$n = p_1^{e_1} \dots p_k^{e_k},$$
>
> where $p_1, \dots, p_k$ are distinct primes and $e_1, \dots, e_k$ are positive integers; this factorisation is unique, apart from permutations of the factors.

**Example 4.4.** $392$ has prime-power factorisation $2^3 7^2$, or alternatively $7^2 2^3$ if we permute the factors, but it has no other prime-power factorisations.

*Proof.* We use **strong induction** to prove the existence of prime-power factorisations. Since we are assuming that $n > 1$, the induction starts with $n = 2$. In this case the factorisation is simply $n = 2^1$. Now assume that $n > 2$ and that every integer strictly between $1$ and $n$ has a prime-power factorisation. If $n$ is prime then there is nothing to do - $n = n^1$ is the required factorisation of $n$. So we can assume that $n$ is composite, which means that $n = ab$ where $1 < a, b < n$. By the induction hypothesis, both $a$ and $b$ have prime-power factorisations, so by substituting these into the equation $n = ab$ and then collecting together powers of each prime $p_i$ we get a prime-power factorisation of $n$. We now prove that the factorisation is in fact unique. Suppose that $n$ has prime-power factorisations

$$n = p_1^{e_1} \dots p_k^{e_k} = q_1^{f_1} \dots q_l^{f_l}$$

where $p_1, \dots, p_k$ and $q_1, \dots, q_l$ are two sets of distinct primes, and the exponents $e_i$ and $f_j$ are all positive integers. The first factorisation shows that $p_1 | n$, so Corollary 4.2 (applied to the second factorisation) implies that $p_1 | q_j$ for some $j = 1, \dots, l$. By renumbering the prime-powers in the second factorisation if necessary, we may assume, for simplicities sake, that $j = 1$. Hence we have that $p_1 | q_1$. Since $q_1$ is prime, it follows that $p_1 = q_1$ and so cancelling this prime from the two factorisations we get

$$p_1^{e_1-1} p_2^{e_2} \dots p_k^{e_k} = q_1^{f_1-1} q_2^{f_2} \dots q_l^{f_l}.$$

We repeat this argument, cancelling primes in the two factorisations, until we run out of primes in one of the factorisations. If one factorisation runs out before the other, then at that stage our reduced factorisations express $1$ as a product of primes $p_i$ or $q_j$, which is impossible since $p_i, q_j > 1$.

It follows that both factorisations run out of primes simultaneously, so we must have cancelled the $e_i$ copies of each $p_i$ with the same number ($f_i$) of

copies of $q_i$. Consequently $k = l$, each $p_i = q_i$, and each $e_i = f_i$ and it follows that the factorisation is unique.  □

**Exercise 4.5.** Find the prime power factorisation of 115, of 188 and of 2020.

As an example of an application of the Fundamental Theorem of Arithmetic, consider the following useful and interesting property (you may have come across the case when $m = 2$).

**Corollary 4.6**

If a positive integer $m$ is not a perfect square then $\sqrt{m}$ is irrational.

*Proof.* It is sufficient to prove **the contrapositive**, that if $\sqrt{m}$ is rational then $m$ is a perfect square. Since we are now considering integers as real numbers we allow ourselves to divide by them (division by $b$ means multiplication by the real number $1/b$) Suppose that $\sqrt{m} = a/b$ where $a$ and $b$ are coprime positive integers and assume, by way of contradiction, that $a, b \geq 2$.

Then the prime factorisations $a = p_1^{e_1} p_2^{e_2} \ldots p_k^{e_k}$, $b = q_1^{f_1} q_2^{f_2} \ldots q_l^{f_l}$ cannot have any primes in common, otherwise the common prime would be a common factor, and coprime integers have no common factors.

Furthermore the prime factorisations of $a^2$ and of $b^2$ can have no common prime factors since they are obtained by squaring the prime factorisations for $a$ and $b$ respectively. Since $mb^2 = a^2$, the prime factorisation of $mb^2$ has the form $mb^2 = p_1^{2e_1} p_2^{2e_2} \ldots p_k^{2e_k}$, so $p_i^{2e_i}$ divides $mb^2$ for each $i$.

But $b^2$ is coprime to $p_i^{2e_i}$ so by Corollary 3.27 (b) $p_i^{2e_i}$ must be a factor of $m$. Furthermore by Corollary 3.27 (a) since the prime factors $p_i^{2e_i}$ are mutually coprime, their product divides $m$, hence $a^2 | m$ and in particular $a^2 \leq m$.

But $b \geq 2$ so $m < 4m \leq mb^2$ and together we get $a^2 \leq m < mb^2 = a^2$. This is a contradiction so our assumption that $\sqrt{m} = a/b$ was false with $a, b \geq 2$.

It follows that if $\sqrt{m} = a/b$ then either $b = 1$ (in which case $\sqrt{m} = a$, and $m = a^2$ as required) or $a = 1$ (in which case $mb^2 = 1$ and so $m = 1 = 1^2$ as required).  □

# Distribution of Primes

> **Theorem 4.7**
>
> There are infinitely many primes.

The fact that there are infinitely many primes, is one of the oldest and most attractive results in mathematics. Here we give the proof given in Book IX of Euclid's Elements.

*Proof.* The proof is by contradiction: we assume that there are only finitely many primes, and then we obtain a contradiction from this, so it follows that there must be infinitely many primes.

Suppose that the only primes are $p_1, p_2, \ldots, p_k$ and let

$$m = p_1 p_2 \ldots p_k + 1.$$

Since $m$ is an integer greater than 1, the Fundamental Theorem of Arithmetic (Theorem 4.3) implies that it is divisible by some prime $p$.

But we assumed that the only primes are $\{p_1, \ldots, p_k\}$ and so $p = p_i$ for some $i$.

Since $p$ divides both $m$ and $p_1 p_2 \ldots p_k$ it divides the difference, $m - p_1 p_2 \ldots p_k = 1$, which is impossible.

Consequently, our initial assumption was false, and so there must be infinitely many primes. $\qquad\square$

**Exercise 4.8.** Let $s_n = n! - 1$. Use the sequence $s_n$ and Theorem 4.3 to show that there are infinitely many primes.

We can use Euclid's proof to obtain a little more information about how frequently prime numbers occur. Let us order the primes with $p_n$ denoting the $n$-th prime, so that $p_1 = 2$, $p_2 = 3$, $p_3 = 5$, and so on.

> **Corollary 4.9**
>
> The $n$-th prime $p_n$ satisfies $p_n \leq 2^{2^{n-1}}$ for all $n \geq 1$.

**NOTE:** This is not a very useful result as the difference between $p_n$ and

$2^{2^{n-1}}$ grows very large, very quickly.

| $n$ | $p_n$ | $2^{2^{n-1}}$ |
|---|---|---|
| 1 | 2 | 2 |
| 2 | 3 | 4 |
| 3 | 5 | 16 |
| 4 | 7 | 256 |
| 5 | 11 | 65536 |
| 6 | 13 | 4294967296 |
| 7 | 17 | 18446744073709551616 |
| 8 | 19 | $3.4028236692093846346337460743177 \times 10^{38}$ |

*Proof.* We use **strong induction** on $n$. The result is true for $n = 1$, since $p_1 = 2 = 2^{2^0}$. Now assume that the result is true for each $n = 1, 2, \ldots, k$.

As in the proof of Theorem 4.7, $p_1 p_2 \ldots p_k + 1$ must be divisible by some prime $p$, which cannot be one of $p_1, p_2, \ldots, p_k$, for then it would divide $1$ - a contradiction.

Now this new prime $p$ must be at least as large as the $(k+1)$-th prime, $p_{k+1}$, so, using the strong induction hypothesis, that $p_i \le 2^{2^{i-1}}$ for $i \le k$,

$$p_{k+1} \le p \le p_1 p_2 \ldots p_k + 1 \le 2^{2^0} 2^{2^1} \ldots 2^{2^{k-1}} + 1$$

Hence:
$$p_{k+1} \le 2^{1+2+4+\cdots+2^{k-1}} + 1.$$

Recalling the formula for the sum of the finite geometric series

$$\sum_{i=0}^{k-1} 2^i = 2^k - 1.$$

we get:
$$p_{k+1} \le 2^{2^k - 1} + 1 = \frac{1}{2} . 2^{2^k} + 1 \le \frac{1}{2} . 2^{2^k} + \frac{1}{2} . 2^{2^k} = 2^{2^k}.$$

This proves the inequality for $n = k + 1$, so by induction it is true for all $n \ge 1$. □

Corollary 4.9 is a *VERY* weak estimate.

By compiling extensive lists of primes, Gauss conjectured in 1793 that the number of primes less than $x$, which we denote by $\pi(x)$, is approximated by the function

$$\text{li} \, x = \int_2^x \frac{dt}{\ln t}$$

or equivalently by $x/\ln x$ in the sense that

$$\frac{\pi(x)}{x/\ln x} \to 1 \quad \text{as} \quad x \to \infty.$$

Here $\ln x = \log_e x$ is the natural logarithm $\int_1^x t^{-1}dt$ of $x$. This result, known as the *Prime Number Theorem*, was eventually proved by Hadamard and de la Vall'ee Poussin in 1896. One can interpret the Prime Number Theorem as showing that the proportion $\pi(x)/\lfloor x \rfloor$ of primes among the positive integers $i \le x$ is approximately $1/\ln x$ for large $x$.

| $x$ | $\pi(x)$ | $\pi(x)/\lfloor x \rfloor$ | $1/\ln(x)$ |
|---:|---:|---:|---:|
| 2 | 1 | 0.5 | 1.44 |
| 3 | 2 | 0.67 | 0.91 |
| 10 | 4 | 0.4 | 0.434 |
| 20 | 8 | 0.4 | 0.334 |
| 50 | 15 | 0.3 | 0.256 |
| 100 | 25 | 0.25 | 0.217 |
| 1000 | 168 | 0.168 | 0.145 |
| 10000 | 1229 | 0.1229 | 0.109 |
| 100000 | 9592 | 0.09592 | 0.087 |

Since $1/\ln(x) \to 0$ as $x \to \infty$, this shows that the primes occur less frequently among larger integers than among smaller integers. For instance there are $168$ primes between $1$ and $1000$, then $135$ primes between $1001$ and $2000$, then $127$ between $2001$ and $3000$, and so on.

**Example 4.10.** Decide how many zeros there are at the end of the integer $100!$ in its decimal notation.

*Solution.* If an integer has exactly $n$ zeros at the end then it is divisible by $10^n$ but not by $10^{n+1}$. It is clear that this is true if it is divisible by $2^n$ and by $5^n$ but not by $2^{n+1}$ and $5^{n+1}$.

Now by the prime factorisation theorem the prime factors of $100!$ are precisely the prime factors of its factors $1, \ldots, 100$. Half of these are even and these each contribute at least one factor of $2$. Hence we obtain $50$ factors of $2$. However of these fifty integers half are divisible by $4$ so these 25 each contribute an additional factor of $2$. Now of these 25, 12 are divisible by $8$ so each contribute three factors of $2$, and of these $6$ contribute four factors of $2$. Three of these contribute five factors of $2$, one contributes $6$, so in total we see $50 + 25 + 12 + 6 + 3 + 1 = 97$ factors of $2$ in $100!$.

Similarly we see $20 + 4 = 24$ factors of $5$ in $100!$. So there are $24$ zeros at the end of $100!$.

# Chapter 5

# Congruences

At the time of typsetting these notes, it was 11 o'clock. In 1 hours time it will be 12 o'clock and in 2 hours it will be 13 o'clock. Of course we don't normally refer to that time as 13 o'clock, but rather as 1 o'clock. In referring to the time, we use a system of arithmetic that we refer to as *modular arithmetic*. In modular arithmetic we simplify number-theoretic problems by replacing each integer with its remainder when divided by some fixed positive integer $n$. So, for example, when telling the time we use $n = 12$.

This has the effect of replacing the infinite number system $\mathbb{Z}$ with a number system $\mathbb{Z}_n$ which contains a finite number ($n$) of elements. In this and the next chapter we shall study the properties of this new number system and show that we can perform arithmetic on these 'numbers' in the same way we can for the integers. We will find that we can

- add,
- subtract and
- multiply

the elements of $\mathbb{Z}_n$, just as we can in $\mathbb{Z}$. Indeed we'll see that the number system $\mathbb{Z}_n$ is an example of a *ring*.

*WARNING: as usual we will avoid "division" but for modular arithmetic there are even some difficulties with cancellation.* The axiom *(ZD)* may not hold.

## 5.1   Congruences

In order to motivate the discussion, let us consider the following two problems.

**Example 5.1.** Show that an integer is divisible by 9 if and only if the sum of its digits is divisible by 9

*Solution.* If a number $x$ is expressed in the form $x_n x_{n-1} \cdots x_0$ this really means the number is equal to $x_n \times 10^n + \cdots + x_1 \times 10 + x_0$. Now $10^n = 9k_n + 1$ for some $k_n$ so $x = (9k_n + 1)x_n + \cdots + (9k_1 + 1)x_1 + x_0$. We can rearrange this as $9(k_n x_n + \cdots + k_1 x_1) + x_n + \cdots + x_0$, so by the remainder theorem $x$ has the same remainder on division by 9 as does $x_n + \cdots + x_0$.

**Example 5.2.** Is $22051946$ a perfect square ?

We could solve this by computing $\sqrt{22051946}$ and determining whether it is an integer, or alternatively by squaring various integers and seeing whether $22051946$ occurs.

There is a much simpler way of seeing that this number cannot be a perfect square.

In Example 3.6 of Chapter 3 we showed that a perfect square must leave a remainder $0$ or $1$ when divided by $4$.

By looking at its last two digits, we see that

$$22051946 = 220519 \times 100 + 46 = 220519 \times 25 \times 4 + 46$$

leaves the same remainder as $46$, and since $46 = 11 \times 4 + 2$ this remainder is $2$. It follows that $22051946$ is *NOT* a perfect square.

**Exercise 5.3.** Show that the last decimal digit of a perfect square cannot be $2, 3, 7$ or $8$. Is $3190491$ a perfect square ?

---

**Definition 5.4**

Let $n$ be a positive integer, and let $a$ and $b$ be any integers. We say that $a$ is *congruent* to $b$ mod $(n)$, or $a$ is a *residue* of $b$ mod $(n)$, and we write

$$a \equiv b \bmod (n),$$

if $a$ and $b$ leave the same remainder when divided by $n$. The number $n$ is referred to as the *modulus* and the arithmetic of congruences, which we shall describe shortly, is called *modular arithmetic*.

---

To calculate residues, we use the Remainder Theorem (Theorem 3.4) to put $a = qn + r$ with $0 \leq r < n$, and $b = q'n + r'$ with $0 \leq r' < n$, and then we say that $a \equiv b \bmod (n)$ if and only if $r = r'$. We will often write simply $a \equiv b$ if it is clear what the value of the modulus $n$ is.

- $100 \equiv 10 \equiv 1 \bmod (9)$ in Example 5.1
- $22051946 \equiv 46 \equiv 2 \bmod (4)$ in Example 5.2.

**Exercise 5.5.** Show that $392 \equiv 1 \bmod (23)$.

As you may expect, we will use the notation $a \not\equiv b \bmod (n)$ to denote that $a$ and $b$ are not congruent mod $(n)$, in other words if they leave different remainders when divided by $n$.

## Casting out nines

Our first example in this chapter can now be reformulated in the following way:

Any integer is congruent to the sum of its digits mod $9$.

We can use this fact to solve the following problem.

**Example 5.6.** Show that rearranging the digits of a power of $2$ cannot yield another power of $2$. (Here we do not allow rearrangements which bring $0$ to the front. )

*Solution.* If we rearrange the digits then the two numbers have the same digit sum and therefore are congruent mod $9$. Now if we examine the table of powers of $2$ we see that $2^0 \equiv 1$, $2^1 \equiv 2$, $2^2 \equiv 4$, $2^3 \equiv 8$, $2^4 \equiv 7$, $2^5 \equiv 5$, $2^6 \equiv 1$ and that this pattern then cycles with period $6$. So if two powers of $2$ are congruent mod $9$ then they are of the form $2^n$ and $2^{n+6k}$ for some $k$. It follows that one of them is at least $64$ times bigger than the other so they cannot have the same number of digits.

## A useful alternative definition of congruence mod $(n)$

**Lemma 5.7**

For any fixed $n \geq 1$ we have $a \equiv b \bmod (n)$ if and only if $n \mid (a - b)$.

*Proof.* Using the Remainder Theorem we let $a = qn + r$ and $b = q'n + r'$ so that

$$a - b = (q - q')n + (r - r')$$

with $-n < r - r' < n$. If $a \equiv b \bmod (n)$ then $r = r'$, so that $r - r' = 0$ and hence $a - b = (q - q')n$, which is divisible by $n$. Conversely, if $n$ divides $a - b$ then it divides the difference $(a - b) - (q - q')n = r - r'$ and since the only integer strictly between $-n$ and $n$ which is divisible by $n$ is $0$, then $r - r' = 0$, which means $r = r'$. Hence $a \equiv b \bmod (n)$ as required. $\qquad\square$

---

**Lemma 5.8**

For any $n \geq 1$,

1. $a \equiv a$ for all integers $a$;
2. if $a \equiv b$ then $b \equiv a$;
3. if $a \equiv b$ and $b \equiv c$ then $a \equiv c$.

---

*Proof.* Let $n \geq 1$.

1. We have $n | (a - a)$ for all $a$.
2. If $n | (a - b)$ then clearly $n | (b - a)$.
3. If $n | (a - b)$ and $n | (b - c)$ then $n | ((a - b) + (b - c)) = a - c$.

$\qquad\square$

## 5.2   Congruence Classes

### Equivalence relations

These three properties are the reflexivity, symmetry and transitivity axioms for an *equivalence relation*. Lemma 5.8 proves that for each fixed $n$, congruence mod $(n)$ is an equivalence relation on $\mathbb{Z}$.

A relation $\sim$ between elements of some set $X$ is said to be an *equivalence relation* if:

- for every $a \in X$, $a \sim a$, (the relation is *reflexive*)
- for every pair $a, b \in X$, if $a \sim b$ then $b \sim a$ (the relation is *symmetric*),
- for every triple $a, b, c \in X$, if $a \sim b$ and $b \sim c$ then $a \sim c$ (the relation is *transitive*).

So the congruence relation $\equiv$ is an equivalence relation on $\mathbb{Z}$. Let $f : X \to Y$ be any function. Consider the relation

$$x \sim y \text{ if and only if } f(x) = f(y).$$

Then $\sim$ is an equivalence relation on $X$.

- $x \sim x$ as $f(x) = f(x)$,

- if $x \sim y$ then $f(x) = f(y)$ and so $f(y) = f(x)$ i.e. $y \sim x$,
- if $x \sim y$ and $y \sim z$ then $f(x) = f(y)$ and $f(y) = f(z)$ and so $f(x) = f(z)$ so that $x \sim z$.

This particular equivalence relation is often referred to as the *kernel* of the function $f$.

# Partitioning a set

The point of an equivalence relation is to give a way of splitting a set up into subsets (pieces). Given any equivalence relation $\sim$ on a set $X$ we define subsets known as the *equivalence classes* as follows:

For each element $x \in X$ let $[x]$ denote the subset $\{y \in X \mid y \sim x\}$

We will show that each equivalence class $[x]$ is a non-empty subset containing $x$ itself and that if $[x] \cap [z] \neq \emptyset$ then $[x] = [z]$. Notice that because $x \sim x$, the equivalence class $[x]$ contains $x$ itself. Now suppose that $[x] \cap [z] \neq \emptyset$. Then there is an element $y \in [x] \cap [z]$. Since $y \in [x]$, $y \sim x$, and since $y \in [z]$, $y \sim z$. Since $\sim$ is symmetric, it follows that $z \sim y$ and $x \sim y$. Since $\sim$ is transitive, $x \sim y \sim z$ implies that $x \sim z$ so $x$ also lies in $[z]$.

Now suppose that $w \in [x]$. Then $w \sim x \sim z$ implies that $w \in [z]$. Since this works for *ANY* $w \in [x]$ we see that $[x] \subseteq [z]$. Running this argument reversing the roles of $x$ and $z$ we can show that $[z] \subseteq [x]$. So $[z] = [x]$ as required.

So have demonstrated that $X = \cup_{x \in X}[x]$ and that $x \in [x]$ for all $x \in X$.

A collection of subsets $P = \{X_i \subseteq X | i \in I\}$ with the property that

1. for each $i, \neq j \in I, X_i \cap X_j = \emptyset$,

2. $X = \cup_{i \in I} X_i$,

is called a *partition* of the set $X$.

**Example 5.9.** Let $X$ be the set of vectors in $\mathbb{R}^2$ and for $u, v \in X$ define $u \sim v$ to mean that $u$ and $v$ have the same length. Then it is clear that

- $u \sim u$ for all $u \in X$;
- if $u \sim v$ then $v \sim u$ for all $u, v \in X$;
- if $u \sim v$ and $v \sim w$ then $u \sim w$, for all $u, v, w \in X$.

Hence $\sim$ is an equivalence relation on $X$. Notice that the equivalence class $[(0, 1)]$ consists of all vectors of length 1 and so is the unit circle in $\mathbb{R}^2$.

<div align="right">

Congruence classes

</div>

Since $\equiv_n$ is an equivalence relation by Lemma 5.8, it follows that $\mathbb{Z}$ is partitioned into disjoint equivalence classes. We refer to these as *congruence classes*.

$$[a] = \{b \in \mathbb{Z} \mid b \equiv a \bmod (n)\}$$

$$= \{\,\ldots, a - 2n, a - n, a, a + n, a + 2n, \ldots\}$$

for $a \in \mathbb{Z}$. Each class corresponds to one of the $n$ possible remainders $r = 0, 1, \ldots, n - 1$ on division by $n$, so there are $n$ different congruence classes. They are

$$[0] = \{\,\ldots, 2n, -n, 0, n, 2n, \ldots\}$$

$$[1] = \{\,\ldots, 1 - 2n, 1 - n, 1, 1 + n, 1 + 2n, \ldots\}$$

$$\vdots$$

$$[n-1] = \{\,\ldots, -n - 1, -1, n - 1, 2n - 1, 3n - 1, \ldots\}$$

There are no further classes distinct from these. For example

$$[n] = \{\,\ldots, -n, 0, n, 2n, 3n, \ldots\} = [0].$$

It should then be clear that the following important observation is then true.

$$[a] = [b] \quad \text{if and only if} \quad a \equiv b \bmod (n)\,.$$

So the following statements all mean the same thing and can be used interchangeably

---

**Proposition 5.10**

The following are equivalent:

- $a \equiv b \bmod (n)$,
- $a$ and $b$ have the same remainder on division by $n$,
- $n \mid (a - b)$,
- $[a] = [b]$.

---

Notice that when $n = 1$ all integers are congruent to each other, so there is a single congruence class, coinciding with $\mathbb{Z}$. So we gain nothing in this case. For $n = 2$ the two classes $[0] = [0]_2$ and $[1] = [1]_2$ consist of the even and odd integers respectively.

## 5.3 Modular Arithmetic

For a given $n \geq 1$, we denote the set of $n$ equivalence classes mod $(n)$ by $\mathbb{Z}_n$. This set is known as the set of integers mod $(n)$. Our next aim is to show how to do arithmetic with these congruence classes, so that $\mathbb{Z}_n$ becomes a number system with properties very similar to those of $\mathbb{Z}$.

We define the binary operations of addition, subtraction and multiplication on the congruence classes in $\mathbb{Z}_n$ in terms of the corresponding operations in $\mathbb{Z}$. If $[a]$ and $[b]$ are elements of $\mathbb{Z}_n$, we define their sum, difference and product to be the classes

- $[a] + [b] = [a + b]$
- $[a][b] = [ab]$

For example, for $n = 3$, $\mathbb{Z}_3 = \{[0], [1], [2]\}$

| $+$ | $[0]$ | $[1]$ | $[2]$ |
|-----|-------|-------|-------|
| $[0]$ | $[0]$ | $[1]$ | $[2]$ |
| $[1]$ | $[1]$ | $[2]$ | $[0]$ |
| $[2]$ | $[2]$ | $[0]$ | $[1]$ |

| $\times$ | $[0]$ | $[1]$ | $[2]$ |
|-----|-------|-------|-------|
| $[0]$ | $[0]$ | $[0]$ | $[0]$ |
| $[1]$ | $[0]$ | $[1]$ | $[2]$ |
| $[2]$ | $[0]$ | $[2]$ | $[1]$ |

### Checking that these operations make sense

Notice that $[2] + [3] = [2]$, that $[5] = [2]$ and $[5] + [3] = [8] = [2]$.

We need to show that these three operations are well-defined, in the sense that the right-hand sides of the three equations defining them depend only on the classes $[a]$ and $[b]$, and not on the particular elements $a$ and $b$ we have chosen from those classes.

In other words, we must show that if $[a] = [a']$ and $[b] = [b']$, then $[a + b] = [a' + b']$, $[a - b] = [a' - b']$ and $[ab] = [a'b']$. These follow immediately from Proposition 5.10 and the following result.

> **Lemma 5.11**
>
> For any $n \geq 1$, if $a' \equiv a$ and $b' \equiv b$ then
>
> - $a' + b' \equiv a + b$
> - $a' - b' \equiv a - b$ and
> - $a'b' \equiv ab$.

*Proof.* If $a' \equiv a$ then $a' = a + kn$ for some integer $k$, and similarly we have $b' = b + ln$ for some integer $l$; then $a' + b' = (a + b) + (k + l)n \equiv a + b$, $a' - b' = (a - b) + (k - l)n \equiv a - b$ and $a'b' = ab + (al + bk + kln)n \equiv ab$.     □

**Exercise 5.12.** Show that each of the axioms [A1] - [A6] for the integers, which we listed in Chapter 1, also hold for $\mathbb{Z}_n$. Show also that axiom [A7] does not, in general, hold.

<div align="right">

What about $[x]^{[y]}$?

</div>

Note that not all arithmetic operations from $\mathbb{Z}$ have well-defined counterparts in $\mathbb{Z}_n$. For example, let us look at what happens if we try to define exponentiation of classes in $\mathbb{Z}_n$ in the obvious way.

We could define

$$[a]^{[b]} = [a^b],$$

restricting $b$ to non-negative values to ensure that $a^b$ is an integer. Notice that $[a]^{[b]}$ is not the same as $[a]^b$ which is equal to $\underbrace{[a][a] \ldots [a]}_{b \text{ times}}$.

So for $n = 3$, we would have, for example,

$$[2]^{[1]} = [2^1] = [2] \, ;$$

**Exponentiation of congruence classes is not well-defined.**

Unfortunately, $[1] = [4]$ in $\mathbb{Z}_3$, and our definition would then mean

$$[2]^{[4]} = [2^4] = [16] = [1] \neq [2] = [2]^{[1]}.$$

Thus we can get different congruence classes for $[a]^{[b]}$ by choosing different elements $b$ and $b'$ in the same class $[b]$. This is because the operation is not well-defined, or in other words, $a' \equiv a$ and $b' \equiv b$ do *NOT* imply $(a')^{b'} \equiv a^b$.

We therefore confine arithmetic in $\mathbb{Z}_n$ to operations which are well-defined, like addition, subtraction, multiplication and integer powers.

We can sometimes cancel or even "divide" in modular arithmetic, but not always so we must be careful.

**Example 5.13.** Let $n = 10$. Then $[7][2] = [14] = [4]$ but also $[2][2] = [4]$ but $[2] \neq [7]$.

So $[7][2] = [2][2]$ but we cannot cancel the factor of $[2]$.

---

**Definition 5.14**

A *multiplicative inverse* for a class $[a] \in \mathbb{Z}_n$ is a class $[b] \in \mathbb{Z}_n$ such that $[a][b] = [1]$. A class $[a] \in \mathbb{Z}_n$ which has a multiplicative inverse is called *unit*.

---

For example for $n = 3$ we saw that $[2][2] = [1]$ so $[2]$ is a unit and is a multiplicative inverse for itself.

Remark / Question: $[n - 1]$ is always a unit and is a multiplicative inverse for itself. Why?

When we have a multiplicative inverse then we **can** cancel a factor.

**Example 5.15.** Let $n = 10$. Then $[3][7] = [21] = [1]$.

So if $[7][a] = [7][b]$ then we can deduce that $[a] = [b]$ because

$$[a] = [1][a] = ([3][7])[a] = [3]([7][a]) = [3]([7][b]) = ([3][7])[b] = [1][b] = [b].$$

One of the main points of Lemma 5.11 is that we can perform rather complicated calculations involving modular arithmetic, one step at a time.

For example to calculate with the expression $12^5 + 89 \times (15^6 - 13^7)$ whilst working modulo 8, instead of expanding the expression to find it equals $-4570603556$ and then trying to calculate which number it is congruence to modulo 8, we can replace different parts of the expression and simplify as we go along.

So $12 \equiv 4 \bmod 8$ and hence $12^5 \equiv 4^5 \bmod 8$. But $4^2 = 16 \equiv 0 \bmod 8$ and so $12^5 \equiv 4^2 \times 4^2 \times 4 \equiv 0 \times 0 \times 4 = 0 \bmod 8$.

In a similar way $89 \equiv 1 \bmod 8$, $15^6 \equiv 1 \bmod 8$, $13^7 \equiv 5 \bmod 8$ and so $12^5 + 89 \times (15^6 - 13^7) \equiv 0 + 1 \times (1 - 5) = -4 \equiv 4 \bmod 8$.

Back to Example

refcnj:04-example-1

Every integer is congruence to the sum of its digits modulo 9. Let $x = 10^n x_n + \ldots + 10 x_1 + x_0$.

Since $10 \equiv 1 (\bmod\ 9)$ then $10^n \equiv 1^n = 1 (\bmod\ 9)$ and so $x \equiv 1 \times x_n \ldots + 1 \times x_1 + x_0 = x_n + \cdots + x_0 (\bmod\ 9)$.

Residues

A set of $n$ integers, containing one representative from each of the $n$ congruence classes in $\mathbb{Z}_n$, is called a *complete set of residues* mod $(n)$. In principle, we can use any complete set of residues but a sensible choice can ease calculations considerably.

The most obvious choice is provided by the Remainder Theorem (Theorem 3.4). We can divide any integer $a$ by $n$ to give $a = qn + r$ for some unique $r$ satisfying $0 \le r < n$; thus each class $[a] \in \mathbb{Z}_n$ contains a unique $r = 0, 1, \ldots, n-1$, so these $n$ integers form a complete set of residues, called the *least non-negative residues* mod $(n)$.

The least absolute residues

For many purposes these are the most convenient residues to use, but sometimes it is better to replace Theorem 3.4 with Exercise 3.5 which gives a remainder $r$ satisfying $-n/2 < r \le n/2$.

These remainders are the *least absolute residues* mod $(n)$, those with least absolute value. When $n$ is odd they are $0, \pm 1, \pm 2, \ldots, \pm(n-1)/2$, and when $n$ is even we also have to include $n/2$. The following calculations illustrate these complete sets of residues.

**Example 5.16.** Let us calculate the least non-negative residue of $26 \times 32$ mod $33$.

Using least absolute residues mod $(33)$, we have $26 \equiv -7$ and $32 \equiv -1$, so Lemma 5.11 implies that $26 \times 32 \equiv (-7) \times (-1) \equiv 7$. Since $0 \le 7 < 33$ it follows that $7$ is the required least non-negative residue.

**Example 5.17.** Calculate the least absolute residue of $15 \times 59$ mod $(75)$.

We have $15 \times 59 \equiv 15 \times (-16)$, and a simple way to evaluate this is to do the multiplication in several stages, reducing the product mod $(75)$ each time. Thus

$$15 \times (-16) = 15 \times (-4) \times 4 = (-60) \times 4 \equiv 15 \times 4 = 60 \equiv -15,$$

and since $-75/2 < -15 \le 75/2$ the required residue is $-15$.

**Example 5.18.** Calculate the least non-negative residue of $7^8$ mod $(17)$.

We do this in several stages, reducing mod $(17)$ whenever possible:

$$7^2 = 49 \equiv -2\,,$$

so that

$$7^4 = (7^2)^2 \equiv (-2)^2 = 4\,,$$

and hence

$$7^8 = (7^4)^2 \equiv 4^2 = 16\,;$$

the required residue is therefore $16$.

**Exercise 5.19.** Without using a calculator

1. find the least absolute residue of $31(33 + 35)$ mod $(23)$;
2. find the least non-negative residue of $31^{33}$ mod $(29)$;
3. explain why the last decimal digit of $1! + 2! + \ldots + n!$ can only take one of 3 possible values and find those values.

# Divisibility and congruences

Since $n$ divides $m$ if and only if $m \equiv 0$ mod $(n)$ if and only if $[m]_n = [0]_n$ it follows that problems about divisibility are equivalent to problems about congruences, and these can sometimes be easier to solve. For example

**Example 5.20.** Let us prove that $a(a + 1)(2a + 1)$ is divisible by $6$ for every integer $a$.

By taking least absolute residues mod $(6)$, we see that if $a$ is any integer, then $a \equiv 0, \pm 1, \pm 2$ or $3$. We examine each possibility. If $a \equiv 0$ then $a(a + 1)(2a + 1) \equiv 0 \times 1 \times 1 \equiv 0$. If $a \equiv 1$ then $a(a + 1)(2a + 1) \equiv 1 \times 2 \times 3 = 6 \equiv 0$. Proceeding in this way we can show that $a(a + 1)(2a + 1) \equiv 0$ in the other four cases, so $6 | a(a + 1)(2a + 1)$ for all $a$.

**Example 5.21.** There is a quicker proof of this, based on the observation that $6 | m$ if and only if $2 | m$ and $3 | m$. We shall make use of this principle later in the course.

*Solution.* For any $a$ either $a$ or $a + 1$ must be even and so congruent to $0$ mod $2$ so $a(a + 1)(2a + 1) \equiv 0$ mod $2$. By the same principle at least one of the three integers $a, a + 1, 2a + 1$ is congruent to $0$ mod $3$. To see this note that $a$ is congruent to $0, 1$ or $2$ mod $3$, so we get $a \equiv 0$ mod $3$, or $2a + 1 \equiv 0$ mod $3$,

or $a + 1 \equiv 0$ mod 3, respectively. Hence the product $a(a+1)(2a+1) \equiv 0$ mod 3, and since 2 and 3 are coprime, $a(a+1)(2a+1) \equiv 0$ mod 6.

The previous argument uses the following more general principle, in which a single congruence mod $(n)$ is replaced with a set of simultaneous congruences mod $(p^e)$ for the various prime powers $p^e$ dividing $n$.

---

**Theorem 5.22**

Let $n$ have prime power factorisation

$$n = p_1^{e_1} \cdots p_k^{e_k},$$

where $p_1, \ldots, p_k$ are distinct primes. Then for any integers $a$ and $b$ we have $a \equiv b$ mod $(n)$ if and only if $a \equiv b$ mod $(p_i^{e_i})$ for each $i = 1, \ldots, k$.

---

We shall delay a formal proof of this until we have met the Chinese Remainder Theorem later in this Chapter. However, as an exercise, you may like to prove it directly, using Corollary 3.27.

Polynomials over $\mathbb{Z}_n$.

---

**Lemma 5.23**

Let $f(x)$ be a polynomial with integer coefficients, and let $n \geq 1$. If $a \equiv b$ mod $(n)$ then $f(a) \equiv f(b)$ mod $(n)$.

---

*Proof.* Write $f(x) = c_0 + c_1 x + \cdots + c_k x^k$, where each $c_i \in \mathbb{Z}$. If $a \equiv b$ mod $(n)$, then repeated use of Lemma 5.11 implies that $a^i \equiv b^i$ for all $i \geq 0$, so $c_i a^i \equiv c_i b^i$ for all $i$, and hence $f(a) = \sum c_i a^i \equiv \sum c_i b^i = f(b)$.  $\square$

For an illustration of this, look at Example 5.20, where we took $f(x) = x(x+1)(2x+1) = 2x^3 + 3x^2 + x$ and $n = 6$; we then used the fact that if $a \equiv 0, \pm 1, \pm 2$ or 3 then $f(a) \equiv f(0), f(\pm 1), f(\pm 2)$ or $f(3)$ respectively, all of which are easily seen to be congruent to 0 mod (6).

Showing that a polynomial has no integer roots

Suppose that a polynomial $f(x)$, with integer coefficients, has an integer root $x = a \in \mathbb{Z}$, so that $f(a) = 0$. It follows then that $f(a) \equiv 0$ mod $(n)$ for all integers $n \geq 1$. The contrapositive of this says:

*if there exists an integer $n \geq 1$ such that the congruence $f(x) \equiv 0$ mod $(n)$ has no solutions, then the equation $f(x) = 0$ has no integer solutions.*

We can often use this to show that certain polynomials $f(x)$ have no integer roots. If $n$ is small we can check whether $f(x) \equiv 0 \bmod (n)$ has any solutions simply by evaluating $f(x_1), \ldots, f(x_n)$ where $x_1, \ldots, x_n$ form a complete set of residues mod $(n)$. Since each $x \in \mathbb{Z}$ is congruent to some $x_i$, then Lemma 5.23 implies that $f(x) \equiv f(x_i)$, and we simply determine whether any of $f(x_1), \ldots, f(x_n)$ is divisible by $n$.

**Example 5.24.** Let us prove that the polynomial $f(x) = x^5 - x^2 + x - 3$ has no integer roots.

To do this, we take $n = 4$ (see later for why we choose 4), and consider the congruence
$$f(x) = x^5 - x^2 + x - 3 \equiv 0 \bmod (4).$$
Using the least absolute residues $0, \pm 1, 2$ as a complete set of residues mod (4), we find that
$$f(0) = -3, \quad f(1) = -2, \quad f(-1) = -6 \quad \text{and} \quad f(2) = 27.$$

None of these values is divisible by 4, so the congruence $f(x) \equiv 0 \bmod (4)$ has no solutions and hence the polynomial $f(x)$ has no integer roots.

Why choose $n = 4$?

The reason is quite simple. For each value of $n < 4$ the congruence $f(x) \equiv 0 \bmod (n)$ *DOES* have a solution $x \in \mathbb{Z}$, even though the equation $f(x) = 0$ does not. So 4 is the smallest value of $n$ for which this method works. If one value of $n$ fails to prove that a polynomial has no integer roots just try a few more values, and if they also fail, this suggests that perhaps there really is an integer root.

**Exercise 5.25.** Prove that the following polynomials have no integer roots:

- $x^3 - x + 1$;
- $x^3 + x^2 - x + 1$;
- $x^3 + x^2 - x + 3$.

**Example 5.26.** Unfortunately, the method used in Example 5.24 is not always strong enough to prove the non-existence of integer roots. For instance, the polynomial
$$f(x) = (x^2 - 13)(x^2 - 17)(x^2 - 221)$$
clearly has no integer roots: indeed, since $13, 17$ and $221 \ (= 13 \times 17)$ are not perfect squares, the roots $\pm\sqrt{13}, \pm\sqrt{17}$ and $\pm\sqrt{221}$ of $f(x)$ are all irrational by Corollary 4.6. However, it can be shown that for every integer $n \geq 1$

there is a solution of $f(x) \equiv 0 \bmod (n)$, so in this case there is no suitable choice of $n$ for our method of congruences.

**Exercise 5.27.** There is no non-constant polynomial $f(x)$, with integer coefficients, such that $f(x)$ is prime for all integers $x$.

## 5.4   Linear congruences

We now return to the question of cancellation of congruence classes, postponed from earlier in this chapter.

We would like to be able to simplify an expression of the form $[a][b] = [a][c]$ in order to conclude that $[b] = [c]$.

This is not always true:

Let $a = 2, b = 3, c = 0$ and $n = 6$. Then $[a][b] = [6] = [0] = [a][c]$, however $[b] = [3] \neq [0] = [c]$. This is reminiscent of the fact that we cannot always cancel even when dealing with integers, since we cannot cancel 0. In the case of mod 6 arithmetic we cannot cancel [2] either.

However we can sometimes cancel, so for example, if $[5][b] = [5][c] \bmod 6$, then $[b] = [c]$. To prove this notice that if $[5][b] = [5][c]$, then $[5][5][b] = [5][5][c]$, so $[25][b] = [25][c]$. But $[25] = [1] \bmod 6$, so $[b] = [c]$.

More generally if we can find a congruence class $[a']$ such that $[a'][a] = [1]$ then we can cancel $[a]$ by multiplying by $[a']$.

<div align="right">Another way to think of this</div>

If $b$ is an integer divisible by an integer $a$, then $b = aq$ for some integer $q$, which we called the quotient of $b$ by $a$. To say that a congruence class $[b]$ is divisible by another congruence class $[a]$ is to say that there is a congruence class $[q]$ such that $[b] = [a][q]$.

<div align="right">Finding divisors</div>

In order to decide whether or not $[b]$ is divisible by $[a]$ we need to consider the solutions of the *linear congruence* $ax \equiv b \bmod (n)$.

Note that if $x$ is a solution, and if $x' \equiv x$, then $ax' \equiv ax \equiv b$ and so $x'$ is also a solution; thus the solutions (if they exist) form a union of congruence classes.

Now $ax \equiv b \bmod (n)$ if and only if $ax - b$ is a multiple of $n$, so $x$ is a solution of this linear congruence if and only if there is integer $y$ such that $x$ and $y$ satisfy the linear diophantine equation $ax + ny = b$, which we studied (with

slightly different notation) in Chapter 3. Translating Theorem 3.31 into the language of congruences we have:

---

**Theorem 5.28**

If $d = \gcd(a, n)$, then the linear congruence

$$ax \equiv b \bmod (n)$$

has a solution if and only if $d$ divides $b$.

If $d$ does divide $b$, and if $x_0$ is any solution, then the general solution is given by

$$x = x_0 + ct$$

where $n = cd$, and $t \in \mathbb{Z}$.

In particular, the solutions form exactly $d$ congruence classes mod $(n)$, with representatives

$$x = x_0,\ x_0 + c,\ x_0 + 2c,\ \ldots,\ x_0 + (d-1)c\,.$$

---

*Proof.* Apart from a slight change of notation, the only part of this which is not a direct translation of Theorem 3.31 is the statement about congruence classes. To prove this, note that

$$x_0 + tc \equiv x_0 + t'c \bmod (n)$$

if and only if $n$ divides $(t - t')c$, that is, if and only if $d$ divides $t - t'$, so the congruence classes of solutions mod $(n)$ are obtained by letting $t$ range over a complete set of residues mod $(d)$, such as $0, 1, \ldots, d - 1$. $\qquad\square$

**Example 5.29.** Solve the congruence

$$10x \equiv 3 \bmod (12)\,.$$

Here $a = 10$, $b = 3$ and $n = 12$, so $d = \gcd(10, 12) = 2$; this does not divide 3, so there are no solutions. (This can be seen directly: the elements of the congruence class $[3]$ in $\mathbb{Z}_{12}$ are all odd, whereas any elements of $[10][x]$ must be even.)

**Example 5.30.** Solve the congruence

$$10x \equiv 6 \bmod (12)\,.$$

As before we have $d = 2$, and now this does divide $b = 6$, so there are two classes of solutions. We can take $x_0 = 3$ as a particular solution, so the general solution has the form

$$x = 3 + \frac{12t}{2} = 3 + 6t \,,$$

where $t \in \mathbb{Z}$. These solutions form two congruence classes $[3]$ and $[9]$ mod $(12)$, with representatives $x_0 = 3$ and $x_0 + (n/d) = 9$; equivalently, they form a single congruence class $[3]$ mod $(6)$.

**Exercise 5.31.** Find the general solution of the congruence $12x \equiv 9$ mod $(15)$.

---

**Corollary 5.32**

If $\gcd(a, n) = 1$ then the solutions of the linear congruence $ax \equiv b$ mod $(n)$ form a single congruence class mod $(n)$.

---

*Proof.* Put $d = 1$ in Theorem 5.28.                                    □

---

**Corollary 5.33**

If $a, n$ are coprime then $[a]$ is a unit in $\mathbb{Z}_n$.

---

*Proof.* From the previous Corollary $ax \equiv 1$ mod $(n)$ has a congruence class $[x]$ as solution. Hence $[a][x] = [1]$.                                    □

## Recap

There are three possibilities:

- If $a$ and $n$ are coprime then for each $b$ there is a unique class $[x]$ such that $[a][x] = [b]$ in $\mathbb{Z}_n$.
- If $d = \gcd(a, n) > 1$ and $d|b$ then there is more than one such class $[x]$
- If $d \nmid b$ there is no such class.

**Example 5.34.** Consider the congruence

$$7x \equiv 3 \text{ mod } (12) \,.$$

Here $a = 7$ and $n = 12$, and since these are coprime there is a single congruence class of solutions; this is the class $[x] = [9]$, since $7 \times 9 = 63 \equiv 3$ mod $(12)$.

<div style="text-align: right;">

Solving $ax \equiv b$ for small $n$

</div>

In Examples 5.29, 5.30 and 5.34, we had $n = 12$. When $n$ is small as in these cases, We can simply calculate $ax$ for each of the $n$ elements $x$ of a complete set of residues mod $(n)$, and see which of these products are congruent to $b$.

<div style="text-align: right;">

Solving $ax \equiv b$ for larger $n$

</div>

In general however, a more efficient method is needed for solving linear congruences. We shall give an algorithm for this, based on Theorem 5.28, but first we need some preliminary results.

---

**Lemma 5.35**

1. Let $m$ divide $a, b$ and $n$, and let $a = a'm$, $b = b'm$ and $n = n'm$. Then

$$ax \equiv b \bmod (n) \quad \text{if and only if} \quad a'x \equiv b' \bmod (n').$$

2. Let $a$ and $n$ be coprime, let $m$ divide $a$ and $b$, and let $a = a'm$ and $b = b'm$. Then

$$ax \equiv b \bmod (n) \quad \text{if and only if} \quad a'x \equiv b' \bmod (n).$$

---

*Proof.*

1. We have $ax \equiv b \bmod (n)$ if and only if $ax - b = qn$ for some integer $q$ and so dividing by $m$, we see that this is equivalent to $a'x - b' = qn'$, that is, to $a'x \equiv b' \bmod (n')$.

2. If $ax \equiv b \bmod (n)$, then as in (1) we have $ax - b = qn$ and hence $(a'x - b')m = qn$. It follows that $m$ divides $qn$ and since $m$ is also coprime to $n$ ($m$ divides $a$ and $a$ is coprime to $n$) then $m$ must divide $q$ by Corollary 3.27(b). Thus, letting $q = q'm$ we have that $a'x - b' = q'n$ is a multiple of $n$, so $a'x \equiv b' \bmod (n)$.
   Conversely, if $a'x \equiv b' \bmod (n)$ then $a'x - b' = q'n$ for some integer $q'$. Hence, multiplying by $m$ we see that $ax - b = mq'n$ and hence $ax \equiv b \bmod (n)$.

$\square$

Note that in (1), where $m$ divides $a, b$ and $n$, we divide all three of these integers by $m$, whereas in (2), where $m$ divides $a$ and $b$, we divide just these two integers by $m$, leaving $n$ unchanged. Also note that in part (2)

$$a'x \equiv b' \bmod n \Rightarrow ax \equiv b \mod n$$

is always true but

$$ax \equiv b \bmod n \Rightarrow a'x \equiv b' \mod n$$

is ONLY true when $\gcd(a, n) = 1$.

**Exercise 5.36.** Show, by means of a counterexample, that Lemma 5.35(2) can fail if $a$ and $n$ are not coprime.

## An algorithm to solve $ax \equiv b \bmod (n)$

To help you understand each step, it may be useful to try this algorithm out on the congruence $10x \equiv 6 \bmod (14)$ - take a look in Appendix E to see if you are correct.

### Step 1

We calculate $d = \gcd(a, n)$, using the Euclidean algorithm if necessary, and see whether $d$ divides $b$. If it does not, there are no solutions, so we stop. If it does, we go on to step 2.

Assuming that there is a solution, Theorem 5.28 will give us the general solution. All we need is to find a particular solution $x_0$, so the rest of the algorithm focusses on finding $x_0$. The general strategy is to reduce $|a|$ until $|a| = 1$, since in this case the solution $x_0 = \pm b$ is obvious.

### Step 2

Since $d$ divides $a, b$ and $n$, Lemma 5.35(1) implies that we can replace the original congruence with

$$a'x \equiv b' \bmod (n'),$$

where $a = a'd$, $b = b'd$ and $n = n'd$. Note also that by Corollary 3.26, $a'$ and $n'$ are coprime. If $d = 1$ then this is nothing practical to do in this step.

### Step 3

We now use Lemma 5.35(ii) to "divide" this new congruence by $m = \gcd(a', b')$, resulting in a congruence

$$a''x \equiv b'' \bmod (n')$$

where $a' = a''m$, $b' = b''m$ and $a''$ is coprime to both $b''$ and $n'$. If $a'' = \pm 1$ then $x_0 = \pm b''$ is the required solution. Otherwise, we go on to step 4. If $m = 1$ then this is nothing practical to do in this step.

Step 4

There are two ways to proceed at this stage.

a. Noting that
$$b'' \equiv b'' \pm n' \equiv b'' \pm 2n' \equiv \ldots \bmod (n') \, ,$$
we may be able to replace $b''$ with some number $b''' \equiv b'' \bmod n'$ such that $\gcd(a'', b''') > 1$.

b. Alternatively, we may be able to multiply the congruence by some suitably chosen constant $c$, giving $ca''x \equiv cb'' \bmod (n')$, in such a way that the least absolute residue $a'''$ of $ca''$ satisfies $|a'''| < |a''|$. Then we have reduced $|a''|$ to give a linear congruence $a'''x \equiv b''' \bmod (n')$ with $b''' = cb''$.

In both cases, we can then apply step 3 to the new congruence $a'''x \equiv b''' \bmod (n')$ and again reduce $|a'''|$. A combination of the methods in step 4 will eventually reduce $a$ to $\pm 1$, in which case the solution $x_0$ can be read off. Theorem 5.28 then gives the general solution.

**Example 5.37.** Consider the congruence

$$10x \equiv 6 \bmod (14) \, .$$

*Solution.* Step 1 gives $\gcd(10, 14) = 2$, which divides 6, so there is a solution. If $x_0$ is any solution, then the general solution is $x = x_0 + (14/2)t = x_0 + 7t$, where $t \in \mathbb{Z}$. To find $x_0$ go onto step 2 - we divide the original congruence through by $\gcd(10, 14) = 2$ to give

$$5x \equiv 3 \bmod (7) \, .$$

Step 3 has no effect since $\gcd(5, 3) = 1$ and so we move on to step 4. We see that $3 \equiv 10 \bmod (7)$, with 10 divisible by 5, we replace the congruence with

$$5x \equiv 10 \bmod (7)$$

and then divide by 5 (which is coprime to 7) to give

$$x \equiv 2 \bmod (7) \, .$$

Thus $x_0 = 2$ is a solution, so the general solution has the form

$$x = 2 + 7t \quad (t \in \mathbb{Z}) \, .$$

**Example 5.38.** Consider the congruence

$$4x \equiv 13 \bmod (47).$$

*Solution.* Step 1 gives $\gcd(4, 47) = 1$, which divides $13$, so the congruence has solutions. If $x_0$ is any solution, then the general solution is $x = x_0 + 47t$ where $t \in \mathbb{Z}$, forming a single congruence class $[x_0]$ in $\mathbb{Z}_{47}$. Since $\gcd(4, 47) = 1$, step 2 has no effect, so we move on to step 3, which also has no effect since $\gcd(4, 13) = 1$. Noting that $4 \times 12 = 48 \equiv 1 \bmod (47)$, we multiply by $12$ to give

$$48x \equiv 12 \times 13 \bmod (47),$$

that is,

$$x \equiv 3 \times 4 \times 13 \equiv 3 \times 52 \equiv 3 \times 5 \equiv 15 \bmod (47).$$

Thus we can take $x_0 = 15$, so the general solution is $x = 15 + 47t$.

For more practice in solving linear congruences, visit Appendix E in the online version of these notes.

## 5.5   Simultaneous linear congruences

In linear algebra, you learn how to solve simultaneous linear equations, for example

$$
\begin{array}{rrrcr}
3x & + & 4y & = & 5 \\
-9x & - & 8y & = & 7
\end{array}
$$

We now consider the solution of simultaneous congruences. In the first century A.D., the Chinese mathematician Sun-Tsu considered problems similar to

*find a number which leaves remainders* $2, 3, 2$ *when divided by* $3, 5, 7$ *respectively.*

In other words, he wanted to find $x$ such that the congruences

$$x \equiv 2 \bmod (3), \quad x \equiv 3 \bmod (5), \quad x \equiv 2 \bmod (7)$$

are simultaneously true. Such problems are thought to have been motivated by observing planetary conjunctions and eclipses.

What do the solutions look like?

Note that if $x_0$ is any solution, then so is $x_0 + (3 \times 5 \times 7)t$ for any integer $t$, so the solutions form a union of congruence classes mod $(105)$.

**Example 5.39.** What are the solutions to the simultaneous congruences

$$x \equiv 3 \bmod (9) , x \equiv 2 \bmod (6)?$$

The simultaneous congruences $x \equiv 3 \bmod (9) , x \equiv 2 \bmod 6$ have no solutions, since if $x \equiv 3 \bmod (9)$ then $3$ divides $x$, whereas if $x \equiv 2 \bmod (6)$ then $3$ does not divide $x$.

## The Chinese Remainder Theorem

The following result, known as the *Chinese Remainder Theorem*, gives a very satisfactory solution to this type of problem.

---

**Theorem 5.40: The Chinese Remainder Theorem**

Let $n_1, n_2, \ldots, n_k$ be positive integers, with $\gcd(n_i, n_j) = 1$ whenever $i \neq j$, and let $a_1, a_2, \ldots, a_k$ be any integers. Then the solutions of the simultaneous congruences

$$x \equiv a_1 \bmod (n_1), \quad x \equiv a_2 \bmod (n_2), \quad \ldots, \quad x \equiv a_k \bmod (n_k)$$

form a single congruence class mod $(n)$, where $n = n_1 n_2 \ldots n_k$.

---

## An application to counting

Suppose we wish to count a large gathering of people, and we know there are around 80 of them. We can do so using the Chinese remainder theorem as follows:

- [Step 1] Get the people to self organize into groups of 7 and count the remainder, $a_1$,
- [Step 2] Now get them to self organize into groups of 13 and count the remainder $a_2$
- [Step 3] If $n$ is the number of people present we now know that $n \equiv a_1 \bmod (7)$ and $n \equiv a_2 \bmod 13$. Since 7 and 13 are coprime we can solve these simultaneous congruences to determine $n \bmod 7.13 = 91$. Since we know that $n$ is around 80 this is enough to determine $n$.

Of course, to carry out the exercise we need to know how to solve simultaneous congruences. The proof of Theorem 3.10 tells us how to do this.

*Proof.* Let $c_i n_i = n$, so $c_i = n_1 \ldots n_{i-1} n_{i+1} \ldots n_k$ for each $i = 1, \ldots, k$. Since each of its factors $n_j$ $(j \neq i)$ is coprime to $n_i$, so is $c_i$. Corollary 5.32 therefore implies that for each $i$, the congruence $c_i x \equiv 1 \bmod (n_i)$ has a single congruence class $[d_i]$ of solutions mod $(n_i)$. We now claim that the integer

$$x_0 = a_1 c_1 d_1 + a_2 c_2 d_2 + \cdots + a_k c_k d_k$$

simultaneously satisfies the given congruences, that is, $x_0 \equiv a_i \bmod (n_i)$ for each $i$. To see this, note that each $c_j$ (other than $c_i$) is divisible by $n_i$, so $a_j c_j d_j \equiv 0 \bmod n_i$ and hence $x_0 \equiv a_i c_i d_i \bmod (n_i)$. But $c_i d_i \equiv 1 \bmod n_i$ and so $x_0 \equiv a_i \bmod n_i$ as required. Thus $x_0$ is a solution of the simultaneous congruences, and it immediately follows that the entire congruence class $[x_0]$ of $x_0 \bmod (n)$ consists of solutions.

To see that this class is unique, suppose that $x$ is any solution. Then $x \equiv a_i \equiv x_0 \bmod (n_i)$ for all $i$ and so each $n_i$ divides $x - x_0$. Since $n_1, \ldots, n_k$ are mutually coprime, repeated use of Corollary 3.27(a) implies that their product $n$ also divides $x - x_0$, so $x \equiv x_0 \bmod (n)$. $\qquad\square$

## Comments

- The proof of Theorem 5.22, which we postponed until later, now follows immediately: given $n = p_1^{e_1} \ldots p_k^{e_k}$, we put $n_i = p_i^{e_i}$ for $i = 1, \ldots, k$, so $n_1, \ldots, n_k$ are mutually coprime with product $n$; the Chinese Remainder Theorem therefore implies that the solutions of the simultaneous congruences $x \equiv b \bmod (n_i)$ form a single congruence class $\bmod (n)$; clearly $b$ is a solution, so these congruences are equivalent to $x \equiv b \bmod (n)$.
- Note that the proof of the Chinese Remainder Theorem does not merely show that there is a solution for the simultaneous congruences; it also gives us a formula for a particular solution $x_0$, and hence for the general solution $x = x_0 + nt$ ($t \in \mathbb{Z}$). We shall see shortly that there is often an easier method to compute the solution.

**Example 5.41.** Consider our original problem

$$x \equiv 2 \bmod (3), \quad x \equiv 3 \bmod (5), \quad x \equiv 2 \bmod (7),$$

*Solution.* We have $n_1 = 3, n_2 = 5$ and $n_3 = 7$, so $n = 105, c_1 = 35, c_2 = 21$ and $c_3 = 15$. We first need to find a solution $x = d_1$ of $c_1 x \equiv 1 \bmod (n_1)$, that is, $35x \equiv 1 \bmod (3)$; this is equivalent to $-x \equiv 1 \bmod (3)$, so we can take $x = d_1 = -1$ for example.

Similarly, $c_2 x \equiv 1 \bmod (n_2)$ gives $21x \equiv 1 \bmod (5)$, that is, $x \equiv 1 \bmod (5)$, so we can take $x = d_2 = 1$, while $c_3 x \equiv 1 \bmod (n_3)$ gives $15x \equiv 1 \bmod (7)$, that is, $x \equiv 1 \bmod (7)$, so we can also take $x = d_3 = 1$.

Of course, different choices of $d_i$ are possible here, leading to different values of $x_0$, but they will all give the same congruence class of solutions $\bmod (105)$.

We now have

$$x_0 = a_1 c_1 d_1 + a_2 c_2 d_2 + a_3 c_3 d_3 = 2.35.(-1) + 3.21.1 + 2.15.1 = 23,$$

so the solutions form the congruence class [23] mod (105), that is, the general solution is $x = 23 + 105t$ ($t \in \mathbb{Z}$).

We can also use the Chinese Remainder Theorem as the basis for a second method for solving simultaneous linear congruences, which is often more efficient.

We start by finding a solution $x = x_1$ to one of the congruences - usually it is best to start with the congruence involving the largest modulus, for reasons that should become apparent. We illustrate by considering the congruences in Example 5.41. Starting with $x \equiv 2$ mod (7), which has $x_1 = 2$ as an obvious solution, and so the 'general solution' for this congruence is

$$x = 2 + 7t, t \in \mathbb{Z}.$$

Notice that there is no real work involved at this point. We want this general solution to also satisfy the other congruences, so taking the next highest modulus we require

$$2 + 7t \equiv 3 \text{ mod } (5),$$

or equivalently $7t \equiv 1$ mod (5). We solve this using our 4-point algorithm as in previous examples: Since $7 \equiv 2$ mod (5) we can simplify to

$$2t \equiv 1 \equiv 6 \text{ mod } (5)$$

and so $t \equiv 3$ mod (5). Hence $t = 3 + 5s$ for $s \in \mathbb{Z}$ and so $x = 2 + 7t = 2 + 7(3 + 5s) = 23 + 35s$ for $s \in \mathbb{Z}$. This is then the general solution for the simultaneous congruences

$$x \equiv 3 \text{ mod } (5), \quad x \equiv 2 \text{ mod } (7).$$

Notice how starting with the modulus with the highest value, reduces the coefficient in front of $t$ immediately, and makes for an easier solution.

We now want this solution to satisfy our final congruence $x \equiv 2$ mod (3) and so we solve

$$23 + 35s \equiv 2 \text{ mod } (3).$$

Since $23 \equiv 2$ mod (3) and $35 \equiv 2$ mod (3) then we have

$$2s \equiv 0 \text{ mod } (3)$$

which has a solution of $s \equiv 0$ mod (3) or $s = 3r$ for $r \in \mathbb{Z}$. Hence $x = 23 + 35(3r) = 23 + 105r$ for $r \in \mathbb{Z}$ is the general solution to all three congruences. Notice that this method only involves solving 2 congruences, whereas the method in the proof of Theorem 5.40 involves solving three (to find the values $d_1, d_2$ and $d_3$).

**Exercise 5.42.** Solve the simultaneous congruences

$$x \equiv 1 \bmod (4), \quad x \equiv 2 \bmod (3), \quad x \equiv 3 \bmod (5).$$

**Exercise 5.43.** Solve the simultaneous congruences

$$x \equiv 2 \bmod (7), \quad x \equiv 7 \bmod (9), \quad x \equiv 3 \bmod (4).$$

## Generalising the CRT

The linear congruences in the Chinese Remainder Theorem are all of the form $x \equiv a_i \bmod (n_i)$. If we are given a set of simultaneous linear congruences, with one (or more) of them in the more general form $ax \equiv b \bmod (n_i)$, then we will first need to use the earlier algorithm to solve this congruence, expressing its general solution as a congruence class modulo some divisor of $n_i$; it will then be possible to apply the techniques based on the Chinese Remainder Theorem to solve the resulting simultaneous congruences.

**Example 5.44.** Consider the simultaneous congruences

$$7x \equiv 3 \bmod (12), \quad 10x \equiv 6 \bmod (14).$$

*Solution.* We saw in Examples 11 and 12 that the first of these congruences has the general solution $x = 9 + 12t$, and that the second has the general solution $x = 2 + 7t$.

It follows that we can replace the original pair of congruences with the pair

$$x \equiv 9 \bmod (12), \quad x \equiv 2 \bmod (7).$$

Clearly, $x_0 = 9$ is a particular solution; since the moduli $12$ and $7$ are coprime, with product $84$, the Chinese Remainder Theorem implies that the general solution has the form $9 + 84t$.

**Exercise 5.45.** Solve the simultaneous congruences

$$3x \equiv 6 \bmod (12), \quad 2x \equiv 5 \bmod (7), \quad 3x \equiv 1 \bmod (5).$$

<div align="right">Simplifying congruences</div>

The Chinese Remainder Theorem can be used to convert a single congruence, with a large modulus, into several simultaneous congruences with smaller moduli, which may be easier to solve.

**Example 5.46.** Consider the linear congruence

$$13x \equiv 71 \bmod (380).$$

*Solution.* Instead of using the algorithm described earlier for solving a single linear congruence, we can use the factorisation $380 = 2^2 \times 5 \times 19$, together with Theorem 5.22, to replace this congruence with the three simultaneous congruences

$$13x \equiv 71 \bmod (4), \quad 13x \equiv 71 \bmod (5), \quad 13x \equiv 71 \bmod (19).$$

These immediately reduce to

$$x \equiv 3 \bmod (4), \quad 3x \equiv 1 \bmod (5), \quad 13x \equiv 14 \bmod (19).$$

The first of these needs no further simplification, but we can apply the single congruence algorithm to simplify each of the other two.

We write the second congruence as $3x \equiv 6 \bmod (5)$, so dividing by $3$ (which is coprime to $5$) we get $x \equiv 2 \bmod (5)$.

Similarly, the third congruence can be written as $-6x \equiv 14 \bmod (19)$, so dividing by $-2$ we get $3x \equiv -7 \equiv 12 \bmod (19)$, and now dividing by $3$ we have $x \equiv 4 \bmod (19)$.

Our original congruence is therefore equivalent to the simultaneous congruences

$$x \equiv 3 \bmod (4), \quad x \equiv 2 \bmod (5), \quad x \equiv 4 \bmod (19).$$

Now these have mutually coprime moduli, so the Chinese Remainder Theorem applies, and we can use either of our two methods to find the general solution.

Using the second method, we start with a solution $x_1 = 4$ of the third congruence; adding and subtracting multiples of $19$, we find that $x_2 = 42$ also satisfies the second congruence, and then adding and subtracting multiples of $19 \times 5 = 95$ we find that $327$ (or equivalently $-53$) also satisfies the first congruence. Thus the general solution has the form $x = 327 + 380t$ ($t \in \mathbb{Z}$).

**Exercise 5.47.** Solve the congruence $91x \equiv 419 \bmod (440)$.

# Chapter 6

# Congruences with a prime modulus

We met an example in the last chapter, where a single congruence mod $(n)$ is equivalent to a set of simultaneous congruences modulo the prime powers $p^e$ appearing in the factorisation of $n$. In this chapter we will study congruences mod $(p)$, where $p$ is prime. In this case we shall see that in addition to modular addition, subtraction and multiplication, cancellation works more smoothly.

## 6.1  Lagrange's Theorem

<div align="right">The arithmetic of $\mathbb{Z}_p$</div>

We saw in Corollary 5.32 that a linear congruence $ax \equiv b \bmod (n)$ has a unique solution mod $(n)$ if $\gcd(a, n) = 1$. Now if $n$ is a prime $p$, then $\gcd(a, n) = \gcd(a, p)$ is either 1 or $p$. In the former case, we have a unique solution mod $(p)$, while in the latter case either every $x$ is a solution (when $p \mid b$) or no $x$ is a solution (when $p \nmid b$).

Looking at this from another point of view, the polynomial $ax - b$ of degree $d = 1$ over $\mathbb{Z}_p$ ($a$ or $b \not\equiv 0 \bmod (p)$), then it has at most one root in $\mathbb{Z}_p$. Now in the Linear Algebra module, we learn that a non-trivial polynomial of degree $d$, with real or complex coefficients, has at most $d$ distinct roots in $\mathbb{R}$ or $\mathbb{C}$. Our first main theorem, due to Lagrange, states that this is also true for the number system $\mathbb{Z}_p$.

> **Theorem 6.1: Lagrange's Theorem**
>
> Let $p$ be prime, and let $f(x) = a_d x^d + \cdots + a_1 x + a_0$ be a polynomial with integer coefficients, where $a_i \not\equiv 0 \bmod (p)$ for some $i$. Then the congruence $f(x) \equiv 0 \bmod (p)$ is satisfied by at most $d$ congruence classes $[x] \in \mathbb{Z}_p$.

1. Note that we allow the possibility that $a_d \equiv 0$, so that $f(x)$ has degree less than $d$.
2. Even if $a_d \not\equiv 0$, $f(x)$ may still have fewer than $d$ roots in $\mathbb{Z}_p$. For example, $f(x) = x^2 + 1$ has only one root in $\mathbb{Z}_2$, namely the class $[1]$, and it has no roots in $\mathbb{Z}_3$.
3. The condition that $a_i \not\equiv 0$ for some $i$ ensures that $f(x)$ yields a nontrivial polynomial when we reduce it mod $(p)$. If $a_i \equiv 0$ for all $i$ then all $p$ classes $[x] \in \mathbb{Z}_p$ satisfy $f(x) \equiv 0$, so the result will fail if $d < p$.
4. If $p$ is not prime the polynomial may have more than $d$ roots. For example, the polynomial $f(x) = x^2 - 1$, of degree $d = 2$, has four roots in $\mathbb{Z}_8$, namely the classes $[1], [3], [5]$ and $[7]$.

*Proof.* We use induction on $d$. If $d = 0$ then $f(x) = a_0$ with $p$ not dividing $a_0$, so there are no solutions of $f(x) \equiv 0$, as required. For the inductive step, we now assume that $d \geq 1$, and that all polynomials $g(x) = b_{d-1} x^{d-1} + \cdots + b_0$ with some $b_i \not\equiv 0$ have at most $d - 1$ roots $[x] \in \mathbb{Z}_p$.

If the congruence $f(x) \equiv 0$ has no solutions, there is nothing left to prove, so suppose that $[a]$ is a solution; thus $f(a) \equiv 0$, so $p$ divides $f(a)$. Now

$$f(x) - f(a) = \sum_{i=0}^{d} a_i x^i - \sum_{i=0}^{d} a_i a^i = \sum_{i=0}^{d} a_i (x^i - a^i) = \sum_{i=1}^{d} a_i (x^i - a^i).$$

For each $i = 2, \ldots, d$ we can put

$$x^i - a^i = (x - a)(x^{i-1} + ax^{i-2} + \cdots + a^{i-2} x + a^{i-1}).$$

By taking out the common factor $x - a$ we have

$$f(x) - f(a) = (x - a)g(x)$$

for some polynomial $g(x)$ with integer coefficients, of degree at most $d - 1$. Now $p$ cannot divide all the coefficients of $g(x)$: if it did, then since it also divides $f(a)$, it would have to divide all the coefficients of $f(x) = f(a) + (x - a)g(x)$, against our assumption. We may therefore apply the induction hypothesis to $g(x)$, so that at most $d - 1$ classes $[x]$ satisfy $g(x) \equiv 0$.

We now count classes $[x]$ satisfying $f(x) \equiv 0$: if any class $[x] = [b]$ satisfies $f(b) \equiv 0$, then $p$ divides both $f(a)$ and $f(b)$, so it divides $f(b) - f(a) = (b - a)g(b)$; since $p$ is prime $p$ divides $b - a$ or $g(b)$, so either $[b] = [a]$ or $g(b) \equiv 0$. There are at most $d - 1$ classes $[b]$ satisfying $g(b) \equiv 0$, and hence at most $1 + (d - 1) = d$ satisfying $f(b) \equiv 0$, as required. $\qquad\square$

**Exercise 6.2.** Find the roots of the polynomial $f(x) = x^2 + 1$ in $\mathbb{Z}_p$ for each prime $p \le 17$. Make a conjecture about how many roots $f(x)$ has in $\mathbb{Z}_p$ for each prime $p$.

We shall revisit this exercise shortly.

---

**Corollary 6.3**

Let $f(x) = a_d x^d + \cdots + a_1 x + a_0$ be a polynomial with integer coefficients, and let $p$ be prime. If $f(x)$ has more than $d$ roots in $\mathbb{Z}_p$, then each coefficient $a_i \equiv 0 \bmod p$.

---

## 6.2 Fermat's Little Theorem

### Polynomials of high degree

Lagrange's Theorem tells us nothing new about polynomials $f(x)$ of degree $d \ge p$: there are only $p$ classes in $\mathbb{Z}_p$, so it is trivial that at most $d$ classes satisfy $f(x) \equiv 0$.

The following result, useful in studying polynomials of high degree, is known as *Fermat's Little Theorem* (not to be confused with Fermat's Last Theorem).

---

**Theorem 6.4: Fermat's Little Theorem**

If $p$ is prime and $a \not\equiv 0 \bmod (p)$, then $a^{p-1} \equiv 1 \bmod (p)$.

---

We give Proof B from (Jones and Jones, 2006), which is purely number-theoretic.

*Proof.* The integers $1, 2, \ldots, p - 1$ form a complete set of non-zero residues mod $(p)$. If $a \not\equiv 0 \bmod (p)$ then $xa \equiv ya$ implies $x \equiv y$, by Corollary 5.32, so that the integers $a, 2a, \ldots, (p-1)a$ lie in distinct classes mod $(p)$. None of these integers is divisible by $p$, so they also form a complete set of non-zero residues.

It follows that $a, 2a, \ldots, (p-1)a$ are congruent to $1, 2, \ldots, p-1$ in some order. (For instance, if $p = 5$ and $a = 3$ then multiplying the residues $1, 2, 3, 4$ by $3$ we get $3, 6, 9, 12$, which are respectively congruent to $3, 1, 4, 2$.) The products of these two sets of integers must therefore lie in the same class, that is,

$$1 \times 2 \times \cdots \times (p-1) \equiv a \times 2a \times \cdots \times (p-1)a \mod (p),$$

or equivalently

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

Rearranging, we get

$$(p-1)!\, (1 - a^{p-1}) \equiv 0 \mod p$$

so since $p$ and $(p-1)!$ are coprime, then we see from Lemma 5.35 (2) that $1 \equiv a^{p-1} \mod p$. □

Theorem 6.4 states that all the classes in $\mathbb{Z}_p$ except $[0]$ are roots of the polynomial $x^{p-1} - 1$.

---

**Corollary 6.5**

If $p$ is prime then $a^p \equiv a \mod (p)$ for every integer $a$.

---

*Proof.* If $a \not\equiv 0$ then Theorem 6.4 gives $a^{p-1} \equiv 1$, so multiplying each side by $a$ gives the result. If $a \equiv 0$ then $a^p \equiv 0$ also, so the result is again true. □

**Example 6.6.** Find the least non-negative residue of $3^{67} \mod (23)$.

Since $23$ is prime and $3$ is not divisible by $23$, we can apply Theorem 6.4 with $p = 23$ and $a = 3$, so that $3^{22} \equiv 1 \mod (23)$. Now $67 = 22 \times 3 + 1$, so

$$3^{67} = (3^{22})^3 \times 3^1 \equiv 1^3 \times 3^1 \equiv 3^1 = 3 \mod (23).$$

**Example 6.7.** Show that $a^{25} - a$ is divisible by $30$ for every integer $a$.

Here Corollary 6.5 is more appropriate, since it refers to all integers $a$, rather than just those coprime to $p$. By factorising $30$, we see that it is sufficient to prove that $a^{25} - a$ is divisible by each of the primes $p = 2, 3$ and $5$. Remember that $p \mid a^{25} - a$ if and only if $a^{25} - a \equiv 0 \mod p$ if and only if $a^{25} \equiv a \mod p$.

Let us deal with $p = 5$ first. Applying Corollary 6.5 twice, we have

$$a^{25} = (a^5)^5 \equiv a^5 \equiv a \mod (5),$$

so $5$ divides $a^{25} - a$ for all $a$. Similarly $a^3 \equiv a \bmod (3)$, so

$$a^{25} = (a^3)^8 a \equiv a^8 a = a^9 = (a^3)^3 \equiv a^3 \equiv a \pmod{3},$$

as required.

For $p = 2$ a direct argument easily shows that $a^{25} - a$ is always even, but we can also continue with this method and use $a^2 \equiv a \bmod (2)$ to deduce (rather laboriously) that

$$a^{25} = (a^2)^{12} a \equiv a^{12} a = (a^2)^6 a = a^6 a =$$

$$(a^2)^3 a \equiv a^3 a = a^4 = (a^2)^2 \equiv a^2 \equiv a \pmod{2}.$$

**Exercise 6.8.** Find the least non-negative residue of $3^{91} \bmod (23)$.

Corollary 6.5 shows that if $f(x)$ is any polynomial of degree $d \geq p$, then by repeatedly replacing any occurrence of $x^p$ with $x$ we can find a polynomial $g(x)$ of degree less than $p$ with the property that $f(x) \equiv g(x)$ for all integers $x$. In other words, when considering polynomials mod $(p)$, it is sufficient to restrict attention to those of degree $d < p$. Similarly, the coefficients can also be simplified by reducing them mod $(p)$.

**Example 6.9.** Let us find all the roots of the congruence

$$f(x) = x^{17} + 6x^{14} + 2x^5 + 1 \equiv 0 \pmod{5}.$$

Here $p = 5$, so by replacing $x^5$ with $x$ we can replace the leading term $x^{17} = (x^5)^3 x^2$ with $x^3 x^2 = x^5$, and hence with $x$. Similarly $x^{14}$ is replaced with $x^2$, and $x^5$ with $x$, so giving the polynomial $x + 6x^2 + 2x + 1$. Reducing the coefficients mod $(5)$ gives $x^2 + 3x + 1$. Thus $f(x) \equiv 0$ is equivalent to the much simpler congruence

$$g(x) = x^2 + 3x + 1 \equiv 0 \pmod{5}.$$

Here we can simply try all five classes $[x] \in \mathbb{Z}_5$, or else note that $g(x) \equiv (x-1)^2$; either way, we find that $[x] = [1]$ is the only root of $g(x) \equiv 0$, so this class is the only root of $f(x) \equiv 0$.

As another application of Fermat's Little Theorem, we prove a result known as **Wilson's Theorem**, though *it was first proved by Lagrange in 1770*:

---

**Corollary 6.10: Wilson's Theorem**

An integer $n > 1$ is prime if and only if $(n-1)! \equiv -1 \bmod (n)$.

---

*Proof.* Suppose that $n$ is a prime $p$. If $p = 2$ then $(p-1)! = 1 \equiv -1 \bmod (p)$, as required, so we may assume that $p$ is odd. Define

$$f(x) = (1-x)(2-x)\ldots(p-1-x) + 1 - x^{p-1},$$

a polynomial with integer coefficients. This has degree $d < p-1$, since when the product is expanded, the two terms in $f(x)$ involving $x^{p-1}$ cancel. If $a = 1, 2, \ldots, p-1$ then $f(a) \equiv 0 \bmod (p)$: the product $(1-a)(2-a)\ldots(p-1-a)$ vanishes since it has a factor equal to $0$, and $1 - a^{p-1} \equiv 0$ by Fermat's Little Theorem.

Thus $f(x)$ has more than $d$ roots mod $(p)$, so by Corollary 6.3 its coefficients are all congruent to $0$ mod $p$. In particular, the constant term $(p-1)! + 1 \equiv 0 \bmod p$, so $(p-1)! \equiv -1$.

For the converse, suppose that $(n-1)! \equiv -1 \bmod (n)$. We then have $(n-1)! \equiv -1 \bmod (m)$ for any factor $m$ of $n$. If $m < n$ then $m$ appears as a factor of $(n-1)!$, so $(n-1)! \equiv 0 \bmod (m)$ and hence $-1 \equiv 0 \bmod (m)$. This implies that $m = 1$, so we conclude that $n$ has no proper factors and is therefore prime. $\qquad\square$

In principle Wilson's theorem gives a method for checking whether or not a number is prime, but in practice it is infeasible for large numbers. The run time of an algorithm based on the theorem (or the storage requirement if it is parallelised) grows very fast with the size of $n$.

In 2003 the Indian mathematician Prof. Manindra Agarwal and two of his students, Nitin Saxena and Neeraj Kayal announced that they had discovered an efficient algorithm for deterministic primality testing. Their paper is only 9 pages long and you can read about it on the web at http://www.cse.iitk.ac.in/users/manindra/algebra/primality_v6.pdf

**Exercise 6.11.** Use Euclid's theorem and Wilson's theorem to evaluate

$$\gcd(29! + 29, 30! + 29).$$

Recall from Exercise 6.2 we saw that the polynomial $x^2 + 1$ in $\mathbb{Z}_p$ has roots:

| $p$ | roots |
|---|---|
| 2 | 1 |
| 3 | |
| 5 | 2 and 3 |
| 7 | |
| 11 | |
| 13 | 5 and 8 |
| 17 | 4 and 13 |

## 6.3 Square roots mod $p$ and mod $pq$

> **Theorem 6.12: The square root of $-1$ mod $p$**
>
> Let $p$ be an odd prime. Then the congruence $x^2 + 1 \equiv 0 \bmod (p)$ has a solution if and only if $p \equiv 1 \bmod (4)$.

*Proof.* Suppose that $p$ is an odd prime, and let $k = (p-1)/2$. In the product

$$(p-1)! = 1 \times 2 \times \cdots \times k \times (k+1) \times \cdots \times (p-2) \times (p-1),$$

we have $p - 1 \equiv -1$, $p - 2 \equiv -2$, $\ldots$, $k + 1 = p - k \equiv -k \bmod (p)$, so by replacing each of the $k$ factors $p - i$ with $-i$ for $i = 1, \ldots, k$ we see that

$$(p-1)! \equiv (-1)^k . (k!)^2 \quad \bmod (p).$$

Now Wilson's Theorem gives $(p-1)! \equiv -1$, so $(-1)^k (k!)^2 \equiv -1$ and hence $(k!)^2 \equiv (-1)^{k+1}$. If $p \equiv 1 \bmod (4)$ then $k$ is even, so $(k!)^2 \equiv -1$ and hence $x = k!$ is a solution of $x^2 + 1 \equiv 0 \bmod (p)$.

On the other hand, suppose that $p \equiv 3 \bmod (4)$, so that $k = (p-1)/2$ is odd. If $x$ is any solution of $x^2 + 1 \equiv 0 \bmod (p)$, then $x$ is coprime to $p$, so Fermat's Little Theorem gives $x^{p-1} \equiv 1 \bmod (p)$. Thus $1 \equiv (x^2)^k \equiv (-1)^k \equiv -1 \bmod (p)$, which is impossible since $p$ is odd, so there can be no solution. $\square$

**Example 6.13.** Let $p = 13$, so $p \equiv 1 \bmod (4)$. Then $k = 6$, and $6! = 720 \equiv 5 \bmod (13)$, so $x = 5$ is a solution of $x^2 + 1 \equiv 0 \bmod (13)$, as is easily verified. The other solution is $-5 \equiv 8 \bmod (13)$.

Lagrange's Theorem implies that if $p$ is any prime then there are at most two classes $[x] \in \mathbb{Z}_p$ of solutions of $x^2 + 1 \equiv 0 \bmod (p)$. When $p \equiv 1 \bmod (4)$ these are the two classes $\pm [k!]$, when $p \equiv 3 \bmod (4)$ there are no solutions, and when $p = 2$ there is a unique class $[1]$ of solutions.

<div align="right">Square roots mod $(p)$ continued</div>

**Proposition 6.14**

Let $p \equiv 3 \bmod (4)$ be prime and let $y$ be an integer with $[y] \neq 0$ in $\mathbb{Z}_p$. Then exactly one of $[y], [-y]$ has square roots in $\mathbb{Z}_p$.

Moreover if $[x] \equiv [y]^{(p+1)/4}$ in $\mathbb{Z}_p$ then $[\pm x]$ are the two square roots of $[y]$ or the two square roots of $[-y]$.

Of course if $[y] = [0]$ in $\mathbb{Z}_p$ then $[0]$ is the unique square root of $[y]$.

*Proof.* As $[y] \neq [0]$ we have $y$ coprime to $p$ so Fermat's theorem tells us that $y^{p-1} \equiv 1 \bmod (p)$. Therefore

$$x^4 \equiv y^{p+1} \equiv y^2 y^{p-1} \equiv y^2 \bmod (p)$$

Hence $x^4 - y^2 \equiv 0 \bmod (p)$. Now factorising we get $(x^2 - y)(x^2 + y) \equiv 0$ so either $x^2 - y \equiv 0 \bmod (p)$ or $x^2 + y \equiv 0 \bmod (p)$: if say $x^2 - y \not\equiv 0 \bmod (p)$ then $x^2 - y$ is coprime to $p$ allowing us to cancel this factor giving $x^2 + y \equiv 0 \bmod (p)$.

Hence $[x]^2 = [-x]^2 = [y]$ or $[x]^2 = [-x]^2 = [-y]$ Lagrange's Theorem tells us that (whichever case we are in) the congruence class mod $(p)$ has at most 2 square roots so $[\pm x]$ are the only square roots of this value.

Now suppose that **both** $[y], [-y]$ have square roots, so $y \equiv a^2$ and $-y \equiv b^2$ for some $a, b \in \mathbb{Z}$.

Let $c = ay^{p-2}$ so $ac = a^2 y^{p-2} = y^{p-1} \equiv 1 \bmod (p)$.

Now

$$(bc)^2 \equiv b^2 c^2 \equiv -yc^2 \equiv -(yc^2) \equiv -(a^2 c^2) \equiv -(ac)^2 \equiv -1 \bmod (p).$$

But by Theorem 6.12, $-1$ has no square roots mod $(p)$ when $p \equiv 3 \bmod (4)$, giving a contradiction.

Hence exactly one of the pair $[y], [-y]$ has a square root mod $(p)$. $\square$

**Example 6.15.** Find the square roots of $5 \bmod (11)$.

Since $11 \equiv 3 \bmod (4)$, we compute $(p+1)/4 = 3$, and set $x \equiv 5^3 \equiv 4 \bmod (11)$. Now $4^2 = 16 \equiv 5 \bmod (11)$ so the square roots of $5$ are $\pm 4 \bmod (11)$.

Now consider the class $2 \bmod (11)$. Carrying out the above calculation we see that if $2$ has any square roots they are $\pm 2^3 \equiv 8 \bmod (11)$. However $8^2 = 64 \equiv 9 \equiv -2 \bmod (11)$, so $2$ has no square roots mod $(11)$.

## Square roots mod $(pq)$

Suppose we want to find the square roots of $71 \bmod (77)$. Factorising $77 = 7 * 11$ we notice that if $x^2 \equiv 71 \bmod (77)$, then $x^2 \equiv 71 \bmod (7)$ and $x^2 \equiv 71 \bmod (11)$, so we have $x^2 \equiv 1 \bmod (7)$ and $x^2 \equiv 5 \bmod (11)$. Applying the proposition we see that $x \equiv \pm 1 \bmod (7)$ and $x \equiv \pm 4 \bmod (11)$.

Now we know from the Chinese remainder theorem that a congruence mod $(7)$ and another mod $(11)$ can be recombined to give us a congruence mod $(77)$, so combining these congruences in the four possible ways we obtain four square roots for $71 \bmod (77)$. For example taking the congruences $x \equiv 1 \bmod (7)$ and $x \equiv 4 \bmod (11)$ we get $x \equiv 15 \bmod (77)$.

Solving all four systems we can recombine to get

$$x \equiv \pm 15, \pm 29 \bmod (77).$$

Hence we have found the four square roots of $71 \bmod (77)$.

## Observations

We solve the 4 pairs of linear congruences

1. $x \equiv 1 \bmod 7, x \equiv 4 \bmod 11$,
2. $x \equiv 1 \bmod 7, x \equiv -4 \bmod 11$,
3. $x \equiv -1 \bmod 7, x \equiv 4 \bmod 11$,
4. $x \equiv -1 \bmod 7, x \equiv -4 \bmod 11$.

But notice that if we put $y = -x$ in (3) then we get $y \equiv 1 \bmod 7, y \equiv -4 \bmod 11$, which we have already solved in (2). So the solution to (3) is just the negative of the solution to (2). Similarly, the solution to (4) is just the negative to the solution for (1).

Notice that $29 \equiv 15 \bmod 7$ and $29 \equiv -15 \bmod 11$, but $29 \not\equiv 15 \bmod 11$ and $29 \not\equiv -15 \bmod 7$. This always happens (providing we are not trying to calculate the square root of a multiple of $p$ or $q$). To see this, suppose that in general, we start with the congruence $x^2 \equiv y \bmod pq$, which then leads to the 4 linear pairs of congruences

1. $x \equiv c \bmod p, x \equiv d \bmod q$, with solution $x \equiv a \bmod pq$,
2. $x \equiv c \bmod p, x \equiv -d \bmod q$, with solution $x \equiv b \bmod pq$,
3. $x \equiv -c \bmod p, x \equiv d \bmod q$, with solution $x \equiv -b \bmod pq$,

4. $x \equiv -c \bmod p, x \equiv -d \bmod q$, with solution $x \equiv -a \bmod pq$.

Now, $a \equiv c \bmod p$ and $b \equiv c \bmod p$ and so $a \equiv b \bmod p$. In a similar way we can see that $a \equiv d \equiv -b \bmod q$ (use (1) & (3)). However, $a \not\equiv -b \bmod p$ since we could then deduce that $b \equiv -b \bmod p$ which would mean that $p|2b$ and so either $p|2$ or $p|b$. But $p$ is an odd prime so can't divide 2 and if $p|b$ then $b \equiv 0 \bmod p$ meaning that $c \equiv b \equiv 0 \bmod p$. But then $y \equiv 0 \bmod p$, a contradiction. In a similar way, $a \not\equiv b \bmod q$.

### Factorising via square roots

Suppose that $n = pq$ is a product of two primes and we know the four solutions $x \equiv \pm a, \pm b$ to the congruence $x^2 \equiv y \bmod (n)$ (assuming that $y$ has a square root mod $n$).

From the previous discussion we know that $a \equiv b \bmod (p)$ and $a \not\equiv b \bmod (q)$ (and $a \equiv -b \bmod (q)$ and $a \not\equiv -b \bmod (p)$) - or the same congruences with $p, q$ interchanged - so that $p|(a-b)$ but $q \nmid (a-b)$. It follows that $\gcd(a-b, n) = p$. Applying the Euclidean algorithm we can therefore extract one of the two factors of $n$.

**Example 6.16.** In the previous example we found the roots $15^2 \equiv 29^2 \equiv 71 \bmod (77)$.

Hence $7 = \gcd(15 - 29, 77)$ is a factor of 77.

### An important principle

All of the calculations we have just done can be done quickly and efficiently (Euclidean algorithm, Chinese remainder Theorem, extracting roots) except factorising $n$. Hence we conclude:

If $n$ is the product of two primes and $y$ is an integer coprime to $n$ which has a square root mod $(n)$, then finding the four square roots of $y$ is computationally equivalent to factorising $n$.

In other words if we can factorise $n$ then we can compute the roots, and if we can compute the roots then we can factorise $n$.

### Square Roots and round coins

We finish the section on square roots with a surprising application.

A story: Alex and Bob are convicts in a high security jail. They have adjacent cells and can talk to one another but cannot see or otherwise contact one another. They both want to read "Escape from Colditz" but the library trolley has only one copy and they can't decide who will get to read it first. I know says Alex "Let's toss a coin!" Alright says Bob fishing a coin out of his pocket and throwing it high into the air. As it lands Alex shouts heads.

Sorry says Bob it was tails. There is a silence, then Alex says "Are you sure?"
. . .

The question is, can we make a fair coin tossing scheme to use in this situation.

The answer is yes, and it uses a technique which relies on the facts we have just proved concerning square roots and factorisation.

## Virtual coin tossing

Alex chooses two large random primes, both congruent to $3$ mod $(4)$. (S)he keeps them secret but tells Bob the product $n = pq$. Bob of course, without access to a large computer has no hope of factorising $n$.

Bob chooses a random integer $x$ and computes $y = x^2$ mod $(n)$. He keeps $x$ secret but tells Alex what $y$ is. Alex, knowing the factorisation of $n$ can compute the four possible roots $\pm a, \pm b$ of $y$. One of these will be $x$ but (s)he doesn't know which.

Alex chooses one of the pairs, either $\pm a$ or $\pm b$ and calls out the result to Bob. Suppose that Alex calls out $\pm a$.

If $x = \pm b$ then Bob tells Alex that Alex has lost. If $x = \pm a$ then Bob tells Alex that Alex has won.

Clearly Alex will win with probability $1/2$ just as in a traditional coin toss.

## But why can't Bob cheat?

Suppose that Alex really loses because $\pm a \not\equiv \pm x$. Then Bob now knows all four roots of $y$, namely $\pm a$ and $\pm x$. Hence he can find a prime factor of $n$ as we discussed above. So Alex can verify that Bob is not cheating by asking him to give him one of the prime factors!

**Example 6.17.** Alex chooses

$$p = 2038074743, q = 1190494759$$

(S)he sends

$$n = pq = 2426317299991771937$$

to Bob. Bob takes

$$x = 141213562373095048$$

and computes

$$y \equiv x^2 \equiv 363278601055491705 \textbf{ mod } (n)$$

which he tells Alex.

Alex computes

$$y^{(p+1)/4} \equiv 1701899961 \bmod (p) \text{ and } y^{(q+1)/4} \equiv 325656728 \bmod (q)$$

The Chinese remainder theorem combines these in four ways to give

$$a \equiv \pm 1012103737618676889, b \equiv \pm 937850352623334103 \bmod (n)$$

Suppose Alex sends 1012103737618676889 to Bob. Then, as this is $-x$ mod $(n)$. Bob declares Alex the winner. If Alex sends 937850352623334103 to Bob then Bob declares Alex the loser and computes

$$\gcd(141213562373095048 - 937850352623334103, n) = 1190494759$$

as proof that he did not cheat!

# Chapter 7

# Euler's Function

In this short section, we consider one of the most important functions in number theory, namely Euler's function $\phi(n)$.

We shall study how to evaluate this function, its basic properties, and see how it can be applied to various problems such as the calculation of large powers and the design and implementation of cipher systems. Although perhaps not immediately obvious, many of the results in Chapter 6 depended on the simple but important fact that if $p$ is prime, and $ab \equiv 0 \bmod (p)$, then $a \equiv 0$ or $b \equiv 0 \bmod (p)$. This makes the arithmetic of $\mathbb{Z}_p$ similar to that of $\mathbb{Z}$, in which the equation $ab = 0$ implies that $a = 0$ or $b = 0$. Unfortunately, this property fails when the modulus is composite. For example, if $n = 12$ then $3 \times 4 \equiv 0 \bmod (12)$ but $3, 4 \not\equiv 0 \bmod (12)$. In Chapter 6 we proved Fermat's Little Theorem, that if $p$ is prime then $a^{p-1} \equiv 1 \bmod (p)$ for all integers $a \not\equiv 0 \bmod (p)$.

It would be useful if we had a similar result for composite moduli. However, if $n = 4$ and $a = 3$ then $a^{n-1} = 27 \not\equiv 1 \bmod (4)$, so the result does not extend to composite values in a straightforward way.

We shall replace $n - 1$ with a different exponent $e(n)$ such that $a^{e(n)} \equiv 1 \bmod (n)$ for all $a$ coprime to $n$.

The simplest function with this property turns out to be Euler's function $\phi(n)$, one of the most important functions in number theory.

## 7.1   Euler's Theorem

Recall that a *multiplicative inverse* for a class $[a] \in \mathbb{Z}_n$ is a class $[b] \in \mathbb{Z}_n$ such that $[a][b] = [1]$.

A class $[a] \in \mathbb{Z}_n$ is a *unit* if it has a multiplicative inverse in $\mathbb{Z}_n$. We let $U_n$ denote the set of units in $\mathbb{Z}_n$.

We sometimes say that the *integer* $a$ is a *unit* mod$(n)$, meaning that $ab \equiv 1$ mod $(n)$ for some integer $b$.

---

**Lemma 7.1**

$[a]$ is a unit in $\mathbb{Z}_n$ if and only if $\gcd(a, n) = 1$.

---

*Proof.*  If $[a]$ is a unit then $ab = 1 + qn$ for some integers $b$ and $q$; any common divisor of $a$ and $n$ would therefore divide $1$, so $\gcd(a, n) = 1$.

Conversely, if $\gcd(a, n) = 1$ then $1 = au + nv$ for some $u$ and $v$ by Theorem 3.16, so $[u]$ is a multiplicative inverse of $[a]$.                              □

**Example 7.2.**

- The units in $\mathbb{Z}_8$ are $[1], [3], [5]$ and $[7]$: in fact $[1][1] = [3][3] = [5][5] = [7][7] = [1]$, so each of these units is its own multiplicative inverse.
- In $\mathbb{Z}_9$, the units are $[1], [2], [4], [5], [7]$ and $[8]$: for instance $[2][5] = [1]$, so $[2]$ and $[5]$ are inverses of each other.

Thus $U_8 = \{[1], [3], [5], [7]\}$ and $U_9 = \{[1], [2], [4], [5], [7], [8]\}$.

### Euler's function

For any positive integer $n$ we define $\phi(n) = |U_n|$, the number of units in $\mathbb{Z}_n$.

By Lemma 7.1 $\phi(n)$ is the number of integers $a = 1, 2, \ldots, n$ such that $\gcd(a, n) = 1$. This function $\phi$ is called *Euler's function*. It is easy to calculate $\phi(n)$ for small $n$. For example,

| $n$ | $=$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\phi(n)$ | $=$ | 1 | 1 | 2 | 2 | 4 | 2 | 6 | 4 | 6 | 4 | 10 | 4 | | |

# Reduced sets of residues

We define a subset $R$ of $\mathbb{Z}$ to be a *reduced set of residues* mod$(n)$ if it contains one element from each of the $\phi(n)$ congruence classes in $U_n$. For instance, $\{1, 3, 5, 7\}$ and $\{\pm 1, \pm 3\}$ are both *reduced* sets of residues mod$(8)$.

If $R$ is a reduced set of residues mod$(n)$, and if an integer $a$ is a unit mod$(n)$, then the set $aR = \{ar \mid r \in R\}$ is also a reduced set of residues mod$(n)$.

To see this note that for any $[r], [s] \in U_n$, $[a][r] = [a][s]$ if and only if $[r] = [s]$, and the inverse of $[a][r]$ is $[b][t]$ where $[b]$ is the inverse of $[a]$ and $[t]$ is the inverse of $[r]$.

> **Theorem 7.3: Euler, 1760**
>
> If $\gcd(a, n) = 1$ then $a^{\phi(n)} \equiv 1$ mod $(n)$.

*Proof.* Proof B of Theorem 6.4 can easily be adapted to this situation; we will merely outline the argument

We replace the integers $1, 2, \ldots, p - 1$ of Theorem 6.4 with a reduced set $R = \{r_1, r_2, \ldots, r_{\phi(n)}\}$ of residues mod$(n)$.

If $\gcd(a, n) = 1$ then $aR$ is also a reduced set of residues mod$(n)$ so the product of all the elements of $aR$ must be congruent to the product of all the elements of $R$. This gives $a^{\phi(n)} r_1 r_2 \ldots r_{\phi(n)} \equiv r_1 r_2 \ldots r_{\phi(n)}$, and since the factors $r_i$ are all units they can be cancelled to give $a^{\phi(n)} \equiv 1$. □

**Example 7.4.** Fermat's Little Theorem is a special case of this result. If $p$ is a prime, then by Lemma 7.1 the units in $\mathbb{Z}_p$ are the classes $[1], [2], \ldots, [p-1]$, so $\phi(p) = p - 1$ and hence $a^{p-1} \equiv 1$ mod $(p)$.

**Example 7.5.** If we take $n = 12$ then $U_{12} = \{\pm[1], \pm[5]\}$, and $\phi(12) = 4$; we have $(\pm 1)^4 = 1$ and $(\pm 5)^4 = 625 \equiv 1$ mod $(12)$, so $a^4 \equiv 1$ mod $(12)$ for each $a$ coprime to $12$.

An interesting application of Euler's Theorem is

> **Theorem 7.6**
>
> If $\gcd(a, n) = 1$ then $a^{\phi(n)-1}$ is an inverse of $a$ modulo $n$.

*Proof.* By Theorem 7.3

$$a^{\phi(n)-1}a = a^{\phi(n)} \equiv 1 \bmod (n)$$

and so $a^{\phi(n)-1}$ is an inverse of $a \bmod (n)$.                     □

**Example 7.7.** The inverse of $5 \bmod (12)$ is $5^3 = 125 \equiv 5 \bmod (12)$.

## 7.2   Euler's Product Formula

### Calculating $\phi(n)$

Computing $\phi(n)$ for small $n$ is straightforward but we need to find a general formula for $\phi(n)$. We have just seen that $\phi(p) = p - 1$ for all primes $p$, and we now extend this to the case where $n$ is a prime-power.

> **Lemma 7.8**
>
> If $n = p^e$ where $p$ is prime, then
>
> $$\phi(n) = p^e - p^{e-1} = p^{e-1}(p-1) = n\left(1 - \frac{1}{p}\right).$$

*Proof.* $\phi(p^e)$ is the number of integers in $\{1, \ldots, p^e\}$ which are coprime to $p^e$, that is, not divisible by $p$; this set has $p^e$ members, of which $p^e/p = p^{e-1}$ are multiples of $p$, so $\phi(p^e) = p^e - p^{e-1} = p^{e-1}(p-1)$.                     □

### Computing $\phi(n)$ for composite $n$

The Fundamental Theorem of Arithmetic, Theorem 4.3, allows us to write any composite number as a product of prime powers. Theorem 7.10 will extend the information given in Lemma 7.8 to give a statement about $\phi(n)$ valid for all integers $n$.

First need the following technical result about complete sets of residues.

> **Lemma 7.9**
>
> If $A$ is a complete set of residues mod $(n)$, and if $m$ and $c$ are integers with $m$ coprime to $n$, then the set $Am + c = \{am + c \mid a \in A\}$ is also a complete set of residues mod $(n)$.

*Proof.* If $am+c \equiv a'm+c \bmod (n)$, where $a, a' \in A$, then by subtracting $c$ and then cancelling the unit $m$, we see that $a \equiv a' \bmod (n)$, and hence $a = a'$. Thus the $n$ elements $am + c$ $(a \in A)$ all lie in different congruence classes, so they form a complete set of residues mod $(n)$. $\qquad\qquad\square$

> **Theorem 7.10**
>
> If $m$ and $n$ are coprime, then $\phi(mn) = \phi(m)\phi(n)$.

*Proof.* We assume that $m, n > 1$, since the result is trivial if $m = 1$ or $n = 1$ ($\phi(1) = 1$). We arrange the $mn$ integers $1, 2, \ldots, mn$ into a rectangular grid with $n$ rows and $m$ columns.

| 1 | 2 | 3 | $\ldots$ | $m$ |
|---|---|---|---|---|
| $m + 1$ | $m + 2$ | $m + 3$ | $\ldots$ | $2m$ |
| $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $(n-1)m + 1$ | $(n-1)m + 2$ | $(n-1)m + 3$ | $\ldots$ | $nm$ |

These integers $1, \ldots, mn$ form a complete set of residues mod $(mn)$, so $\phi(mn)$ is the number that are coprime to $mn$. But this amounts to counting the number of these residues coprime to both $m$ and $n$.

### Counting the entries mod $m$

The integers in a given column are all congruent mod $(m)$, and the $m$ columns correspond to the $m$ congruence classes mod $(m)$. Hence exactly $\phi(m)$ of the columns consist of integers $i$ coprime to $m$, and the other columns consist of integers with $\gcd(i, m) > 1$.

### Counting them mod $n$

Now each column of integers coprime to $m$ has the form $c, m + c, 2m + c, \ldots, (n-1)m+c$ for some $c$. By Lemma 7.9 this is a complete set of residues mod $(n)$, since $A = \{0, 1, 2, \ldots, n - 1\}$ is and since $\gcd(m, n) = 1$. Such a column therefore contains $\phi(n)$ integers coprime to $n$, then these $\phi(m)$

columns yield $\phi(m)\phi(n)$ integers $i$ coprime to both $m$ and $n$. Thus $\phi(mn) = \phi(m)\phi(n)$, as required.                                                                                   □

**Example 7.11.** The integers $m = 3$ and $n = 4$ are coprime, with $\phi(3) = \phi(4) = 2$; here $mn = 12$ and $\phi(12) = 2.2 = 4$.

**Exercise 7.12.** Form the array in the above proof with $m = 5$ and $n = 4$; by finding the entries coprime to 20, verify that $\phi(20) = \phi(5)\phi(4)$.

---

**Corollary 7.13**

If $n$ has prime-power factorisation $n = p_1^{e_1} \ldots p_k^{e_k}$ then

$$\phi(n) = \prod_{i=1}^{k} \phi(p_i^{e_i}) = \prod_{i=1}^{k} (p_i^{e_i} - p_i^{e_i-1}) = n \prod_{i=1}^{k} \left(1 - \frac{1}{p_i}\right).$$

---

*Proof.* We prove the first expression by induction on $k$, the other expressions then follow easily.

For the anchoring step with $k = 1$ we use Lemma 7.8. Now assume that $k > 1$ and that the result is true for all integers divisible by fewer than $k$ primes.

We have $n = p_1^{e_1} \ldots p_{k-1}^{e_{k-1}} . p_k^{e_k}$, where $p_1^{e_1} \ldots p_{k-1}^{e_{k-1}}$ and $p_k^{e_k}$ are coprime, so Theorem 7.10 gives

$$\phi(n) = \phi(p_1^{e_1} \ldots p_{k-1}^{e_{k-1}})\phi(p_k^{e_k}).$$

The induction hypothesis gives

$$\phi(p_1^{e_1} \ldots p_{k-1}^{e_{k-1}}) = \prod_{i=1}^{k-1} (p_i^{e_i} - p_i^{e_i-1}),$$

Now Lemma 7.8 gives

$$\phi(p_k^{e_k}) = (p_k^{e_k} - p_k^{e_k-1}),$$

so by combining these two results we get

$$\phi(n) = \prod_{i=1}^{k} (p_i^{e_i} - p_i^{e_i-1}).$$

□

### An alternative form of the result

We often express this result more concisely using the notation $\phi(n) = n \prod_{p|n}(1 - \frac{1}{p})$, where $\prod_{p|n}$ denotes the product over all primes $p$ dividing $n$.

**Example 7.14.** Calculate $\phi(60)$.

The primes dividing $60$ are $2, 3$ and $5$, so

$$\phi(60) = 60\left(1 - \frac{1}{2}\right)\left(1 - \frac{1}{3}\right)\left(1 - \frac{1}{5}\right) = 60.\frac{1}{2}.\frac{2}{3}.\frac{4}{5} = 16\,.$$

We can confirm this by writing down the integers $i = 1, 2, \ldots, 60$, and then deleting those with $\gcd(i, 60) > 1$.

**Exercise 7.15.** Calculate $\phi(42)$, and confirm it by finding a reduced set of residues $\mod(42)$.

**Exercise 7.16.** For which values of $n$ is $\phi(n)$ odd? Show that there are integers $n$ with $\phi(n) = 2, 4, 6, 8, 10$ and $12$, but not $14$.

**Exercise 7.17.** Show that for each integer $m$, there are only finitely many integers $n$ such that $\phi(n) = m$.

**Exercise 7.18.** Find the smallest integer $n$ such that $\phi(n)/n < 1/4$.

### Applications of Euler's function

Having seen how to calculate Euler's function $\phi(n)$, we now look for some applications of it. We saw in Chapter 6 how to use Fermat's Little Theorem $a^{p-1} \equiv 1$ to simplify congruences $\mod(p)$, where $p$ is prime, and we can now make similar use of Euler's Theorem $a^{\phi(n)} \equiv 1$ to simplify congruences $\mod(n)$ when $n$ is composite.

**Example 7.19.** Find the last two decimal digits of $3^{1492}$.

This is equivalent to finding the least non-negative residue of $3^{1492} \mod (100)$.

Now $3$ is coprime to $100$, so Theorem 7.3 (with $a = 3$ and $n = 100$) gives $3^{\phi(100)} \equiv 1 \mod (100)$. The primes dividing $100$ are $2$ and $5$, so Corollary 7.13 gives $\phi(100) = 100.(1/2).(4/5) = 40$, and hence we have $3^{40} \equiv 1 \mod (100)$. Since $1492 \equiv 12 \mod (40)$, it follows that $3^{1492} \equiv 3^{12} \mod (100)$. Now $3^4 = 81 \equiv -19 \mod (100)$, so $3^8 \equiv (-19)^2 = 361 \equiv -39$ and hence $3^{12} \equiv -19. - 39 = 741 \equiv 41$. The last two digits are therefore $41$.

**Exercise 7.20.** Show that if a positive integer $a$ is coprime to $10$, then the last three decimal digits of $a^{2001}$ are the same as those of $a$.

# Chapter 8

# Applications to cryptography

Secret ciphers have been used since ancient times to send messages securely, for instance in times of war or diplomatic tension.



Nowadays sensitive information of a medical or financial nature is often stored in computers, and it is important to keep it secret.

Throughout this section, we shall assume that the alphabet that we are working with only consists of the upper case letters A,B,...,Z.

## Caesar shift ciphers

Many ciphers are based on number theory. A simple one is to replace each letter of the alphabet with its successor. Mathematically, we can do this by representing

the letters as integers, say A=0, B=1, ..., Z=25, and then adding $1$ to each. In order to encipher Z as A, we must add mod $(26)$, so that $25 + 1 \equiv 0$. Similar ciphers are obtained by adding some fixed integer $k$ (known as the *key*), rather than $1$: Julius Caesar used the key $k = 3$. To decipher, we simply apply the reverse transformation, subtracting $k$ mod $(26)$.

ATTACK TOMORROW $\rightarrow$ DWWDFN WRPRUURZ

Figure 8.1: Julius Caesar

## Frequency Analysis

There are two problems with the Caesar Shift cipher. The first is that the *key space* is small. There are only 25 non-trivial values for the key and we can easily exhaust all values in order to test which one was actually used.

The second problem concerns *frequency analysis*. In an 'average' piece of English language text, the letters tend to appear with a definite frequency. The letter 'E' is usually the most frequent, followed by the letter 'T', then 'A' and so on.



If we shift each character by a set amount, then the shifted value of 'E' is likely to be the most frequent, followed by the shifted value of the letter 'T' etc.

We only need to have an educated guess as to what one character has been shifted by, to discover the key used.

### Affine shift ciphers

A slightly more secure class of ciphers uses affine transformations of the form $x \mapsto ax + b \bmod (26)$, for various integers $a$ and $b$.

To decipher successfully, we need to be able to recover the value of $x$ uniquely from $ax + b$. More precisely we need to know that if $ax + b \equiv ay + b \bmod 26$ then $x \equiv y \bmod 26$.

In the language of functions we want the affine transformation $f : \mathbb{Z}_{26} \longrightarrow \mathbb{Z}_{26}$ defined by the rule $x \mapsto ax + b$ to be one to one, or *injective*.

If we take $a = 5, b = 3$ we have

ATTACK TOMORROW $\rightarrow$ DUUDNB UVLVKKVJ

### When is an affine shift injective?

Suppose that $ax + b \equiv ay + b \bmod 26$. Then, subtracting $b$ from both sides we see that $ax \equiv ay \bmod 26$, and so $a(x - y) \equiv 0 \bmod 26$. This equation has the obvious solution $x = y$, and if $a$ is coprime to 26, then this is the only solution, since we may cancel the unit $a$. If on the other hand $a$ is not a unit, then $\gcd(a, 26) = d > 1$ and we may choose $r$ such that $rd = 26$ and $s$ such that $a = sd$. Now it is clear that $ra = rds = 26s \equiv 0 \bmod 26$, and so for any $x$ there is a $y = x - r$ such that $a(x - y) = ar \equiv 0 \bmod 26$, yielding a second solution to the equation. We conclude that the affine shift is injective if and only if $a$ is a unit mod $(26)$.

By counting the pairs $a, b$ we see that there are $\phi(26) \times 26 = 12 \times 26 = 312$ such ciphers. Breaking such a cipher by trying all the possibilities for $a$ and $b$ would be tedious by hand (though simple with a computer), but again frequency searches can make the task much easier.

**Exercise 8.1.** If the enciphering transformation is $x \mapsto 7x + 3 \bmod (26)$, encipher GAUSS and decipher MFSJDG.

## 8.1 Ciphers using modular exponentiation

### Discrete log encryption

We can do rather better with ciphers based on Fermat's Little Theorem.

The idea is as follows. We choose a large prime $p$, and an integer $e$ coprime to $p - 1$. For encoding, we use the transformation $\mathbb{Z}_p \to \mathbb{Z}_p$ given by $x \mapsto x^e$ mod $(p)$.

(We saw in Chapter 6 how to calculate large powers efficiently in $\mathbb{Z}_p$.)

With $p = 29, e = 5$

ATTACK TOMORROW $\to$ AVVADI VTMTRRTN

## Deciphering the message

If $0 < x < p$ then $x$ will be coprime to $p$, so $x^{p-1} \equiv 1 \mod (p)$. To decipher, we first find the multiplicative inverse $f$ of $e$ mod $(p - 1)$, that is, we solve the congruence $ef \equiv 1 \mod (p - 1)$, using the method described in Chapter 6; this is possible since $e$ is a unit mod $(p - 1)$. Then $ef = (p - 1)k + 1$ for some integer $k$, so

$$(x^e)^f = x^{(p-1)k+1} = (x^{p-1})^k . x \equiv x \text{ deciphered efficiently}$$

.

**Example 8.2.** Suppose that $p = 29$ (unrealistically small, but useful for a simple illustration). We must choose $e$ coprime to $p - 1 = 28$, and then find $f$ such that $ef \equiv 1 \mod (28)$. If we take $e = 5$, for example, so that encoding is given by $x \mapsto x^5 \mod (29)$, then $f = 17$ and decoding is given by $x \mapsto x^{17} \mod (29)$. Note that

$$(x^5)^{17} = x^{85} = (x^{28})^3 . x \equiv x \mod (29)$$

since $x^{28} \equiv 1 \mod (29)$ for all $x$ coprime to $29$, so decoding is the inverse of encoding.

**Exercise 8.3.** In Example 8.2, encipher $9$ and decipher $11$.

## Block ciphers

Representing individual letters as numbers tends to be insecure, since an eavesdropper could use known frequencies of letters. A better method is to group the letters into blocks of length $k$, and to represent each block as an integer $x$. (If the length of the message is not divisible by $k$, one can always add extra meaningless letters at the end.)

We choose $p$ sufficiently large that the distinct blocks of length $k$ can be represented by different congruence classes $x \not\equiv 0 \mod (p)$, and then the encoding and decoding are given as before by $x \mapsto x^e \mod (p)$ and $x \mapsto x^f \mod (p)$.

Example: We encode pairs of letters as follows

$$\text{AA}\mapsto 1, \text{AB}\mapsto 2, \ldots, \text{AZ}\mapsto 26, BA \mapsto 27, \ldots, \text{ZZ}\mapsto 26^2 = 676.$$

With $p = 677, e = 7$ and $k = 2$

ATTACK TOMORROW $\rightarrow$ GFYOMA RHYDVRMH Breaking this cipher seems to be very difficult. Suppose, for instance, that an eavesdropper has discovered the value of $p$ being used, and also knows one pair $x$ and $y \equiv x^e \bmod (p)$.

To break the cipher, the eavesdropper needs to know the value of $f$ (or equivalently $e$).

If $p$ is sufficiently large (say a few hundred decimal digits) then there is no known efficient algorithm for calculating $e$ from the congruence $y \equiv x^e \bmod (p)$, where $x, y$ and $p$ are known. This is sometimes called the *discrete logarithm problem*, since we can regard this congruence as a modular version of the equation $e = \log_x(y)$.

The whole point of this cipher is that, while exponentials are easy to calculate in modular arithmetic, logarithms are *apparently* difficult.

**Exercise 8.4.** Find a value of $e$ coprime to $28$ such that $27 \equiv 10^e \bmod (29)$.

# Key exchange

The one weakness of this type of cipher is that the sender and receiver must first agree on the values of $p$ and $e$ (called the *KEY* of the cipher) before they can use it. How can they do this secretly, bearing in mind that they will probably need to change the key from time to time for security ? They could, of course exchange this information in enciphered form, but then they would have to agree about the details of the cipher used for discussing the key, so they are no nearer solving the problem. Alice and Bob agree to use a prime number $p$ and base $g$. Alice chooses a secret number $a$ and sends Bob the number

$$A = g^a \mod p.$$

Bob chooses a secret number $b$ and sends Alice

$$B = g^b \mod p.$$

Alice calculates $B^a \mod p$ and Bob calculates $A^b \mod p$. This is their secret key.

For example, if $p = 29, g = 5$. If Alice chooses $a = 7$ and Bob chooses $b = 11$, then

$$A = 5^7 \mod 29 = 28$$

while

$$B = 5^{11} \mod 29 = 13.$$

Now $s = 28^{11} \mod 29 = 28 = 13^7 \mod 29$. This is referred to as the *Diffie-Helman Key Exchange protocol* or sometimes the *Diffie-Helman-Merkle Key Exchange protocol*.

## Public key cryptography

One can avoid this difficulty by using a *public-key cryptographic system*. Each person using the system publishes numerical information which enables any other user to encipher messages, without giving away sufficient information to allow anyone but him/herself to decipher them.

The most famous of these systems is the RSA encryption system publicly described first by Rivest, Shamir and Adelman of MIT in 1977. An equivalent system had been described in a top secret document at GCHQ by Clifford Cocks in 1973.

## Symmetric v Asymmetric encryption schemes

An encryption method is said to be *symmetric* if anyone in possession of the encryption algorithm and key is in a position to decrypt all messages encrypted with it.

Examples:

- Affine shift ciphers - given the algorithm $x \to ax + b \mod n$ and the constants $a, b, n$ one can compute the inverse function $x \to c(x - b)$ by solving the equation $ac \equiv 1 \mod n$.
- Discrete log ciphers - given the algorithm $x \to x^e \mod p$ and the constants $e, p$ one can compute the inverse function $x \to x^f$ by solving the equation $ef \equiv 1 \mod p - 1$.

## The Enigma Machine

Notice that some historical ciphers were symmetric in a stronger sense, in that the encryption and decryption functions were the same function! Simple examples are given by the shift cipher $x \to x + 13$, and the affine shift $x \to -x$. A much more sophisticated example was given by the enigma machine. However it was setup to encrypt a message, the same setting (or **key**) would be used to decrypt it.

## RSA encryption

As in the discrete log system the encryption is achieved by taking powers.

To set up an RSA cipher one chooses a distinct pair of large primes $p$ and $q$, calculates $n = pq$, and chooses a positive integer $e$ coprime to $\phi(n) = (p-1)(q-1)$.

The encryption algorithm is given by $x \to x^e \mod (n)$ and the key, which is freely published is the pair $(e, n)$.

Example: The encryption key might be $(7, 143)$, but it could not be $(3, 143)$ or $(5, 143)$. Why not?

## Deciphering RSA

As with the discrete log cipher the decryption algorithm is given by $x \to x^f \bmod (n)$ for a suitable choice of exponent $f$.

This choice must satisfy $(x^e)^f \equiv x \bmod (n)$ for every $x$, and it should be clear that Euler's theorem tells us how to choose $f$.

We know that for $x$ coprime to $n$, $x^{t.\phi(n)+1} \equiv (x^{\phi(n)})^t.x \equiv 1^t.x \equiv x \bmod (n)$, so at least for these values of $x$ it is sufficient to choose $f$ so that $ef \equiv 1 \bmod (\phi(n))$.

In fact, whether or not $x$ is coprime to $n$, $x^{ef} \equiv x^{t(p-1)(q-1)+1} \equiv x \mod p$ and $x^{ef} \equiv x^{t(p-1)(q-1)+1} \equiv x \mod q$. So we see that both $p$ and $q$ divide the difference $x^{ef} - x$, and since they are coprime Corollary 3.27 tells us that their product also divides this difference so $x^{ef} \equiv x \mod pq$ as required, for all $x$.

Now if we set up the RSA cipher then we know that $n = pq$ so we can compute $\phi(n) = (p-1)(q-1)$ and can therefore solve the equation for $f$ in terms of $e, p, q$. As we will next see this is essentially the only way we can compute $f$!

**Why we (think we) know that RSA decryption is hard for the enemy**

Suppose the enemy is able to compute $\phi(n)$ and therefore to solve the equation $ef \equiv 1 \mod \phi(n)$ thereby finding the decryption key.

Then the enemy knows $n$ and $\phi(n)$. They therefore know the integers $pq$ and $(p-1)(q-1)$. Taking the difference and adding $1$ gives $p + q$, so the enemy knows the two integers $pq$ and $p + q$. This is enough to compute the prime factorisation of $n$.

To see this note that

$$(p - q)^2 = (p + q)^2 - 4pq.$$

From this we may extract $p - q$ as well. Then knowing $p + q$ and $p - q$ we can compute both $p$ and $q$ as required.

## The prime factorisation problem

Given an integer $n$, how can we find the prime factors of $n$?

The prime factorisation problem is hard. The best methods currently take around $18$ months of parallel processing (around $50$ years of processor

time) to factor an integer with around $200$ digits. Typical $1024$ bit RSA exponents have around $300$ digits. In around 2007 Lenstra and his collaborators cracked a 200 digit number. Asked whether 1024 bit RSA ciphers were dead he replied:

**"The answer to that question is an unqualified yes.""**

**Example 8.5.** Suppose that $p = 89$ and $q = 97$ are chosen, so $n = 89 \times 97 = 8633$ is published, while $\phi(n) = 88 \times 96 = 8448 = 2^8 \times 3 \times 11$ is kept secret. The receiver chooses and publishes an integer $e$ coprime to $\phi(n)$, say $e = 71$. He then finds (and keeps secret) the multiplicative inverse $f = 119$ of $71 \bmod (8448)$.

This can be done using the Euclidean algorithm.

To check the answer, note that $71 \times 119 = 8449 \equiv 1 \bmod (8448)$. To send a message, anyone can look up the pair $n = 8633, e = 71$, and use the encoding $x \mapsto x^{71} \bmod (8633)$. The receiver uses the decoding transformation $x \mapsto x^{119} \bmod (8633)$, which is not available to anyone who does not know that $f = 119$. An eavesdropper would need to factorise $n = 8633$ in order to find $\phi(n)$ and then $f$. Of course, factorising $8633$ is not so difficult, but this is just a simple illustration of the method, and significantly larger primes $p$ and $q$ would pose a much harder problem.

**Factorising** $8633$, **using the fact that** $\phi(8633) = 8448$

As in the description above we know $pq = 8633$ and $(p-1)(q-1) = 8448$, so $pq - (p+q) + 1 = 8448$. From this we conclude that $8633 - 8448 + 1 = p + q$, so we put $p + q = 186$.

Now we see that $(p-q)^2 = (p+q)^2 - 4pq = 186^2 - 4.8633 = 64$ so $p - q = \pm 8$. Now $2p = (p+q) + (p-q) = 186 + 8 = 194$ showing us that $p = 97$. Similarly $2q = (p+q) - (p-q) = 186 - 8 = 178$ so $q = 89$.

Of course our calculation has reversed the primes but this does not matter.

**Exercise 8.6.** If my public key is the pair $n = 10147, e = 119$, then what is my decoding transformation ?

# 8.2   Digital Signatures

## Cryptographic signatures/digital certificates

This system also gives a way of *signing* a message, to prove to a receiver that it comes from you and from nobody else.

The fundamental problem with digital signatures is the question of how you stop them being cut and paste from one document you did send onto another one you didn't. Once this problem is solved (which is done using hash functions, described below) the idea is to turn public key encryption on its head, publishing a decryption key that everyone knows and keeping secret an encryption key.

## Digests

The first stage in digitally signing the document is to produce a so called **digest** (also known as a hash), that is, a short summary of the document, typically 1024 bits long. The method of producing the digest varies between the different secure signature systems, but it is an algorithm chosen to make it unlikely that two different messages will produce the same digest. This digest is attached to the message and **proves** that it has not been tampered with.

So long as the message and its digest match we have reason to believe the message is authentic.

But the digest algorithm is public so anyone tampering with the message just has to also replace the digest too. To prevent that, the digest is encrypted using the secret encryption algorithm.

## Verifying the message

On receipt of the message the recipient deciphers the digest using the published decryption key, and then uses the digest algorithm to produce their own version of the message digest directly from the message contents. If the decrypted digest matches the new one then the signature and the message (almost certainly) belong together and the message is genuine. The security of these systems relies in part on the reliability of the digest used.

## The Birthday Attack

This is a method to circumvent digital signatures based on the Birthday paradox. The classical birthday paradox states that in a group of 23 people there is a probability of 1/2 that at least two of them share the same birthday. In a larger group the odds are even higher.

Assuming that birth dates are uniformly distributed (actually they are not) and ignoring leap years there are 365 possible birthday dates. Consider all the possible collections of birthdays assuming that no two people in a group of size $n$ share a birthday.

The first person in the group has 365 possible birthdays, the second has 364 and so on. Hence there are $365!/(365 - n)!$ possibilities for the list of birthdays, so the probability that no two people share a birthday is:

$$\frac{365!/(365 - n)!}{365^n}.$$

It follows that the probability that at least two people have a shared birthday is

$$1 - \frac{365!/(365 - n)!}{365^n}.$$



## Attacking the digital signature

The security of the digital signature depends on the assumption that different messages have different hashes. If you have a large collection of messages then, as with birthdays, the odds that two messages have the same hash is higher than you might think. The birthday attack exploits this weakness.

If a 64 bit hash is used, that is each message is converted to a digest which is stored as a binary number with 64 digits, there are $2^{64}$ (or approximately $1.8 \times 10^{19}$) different digests. If these are all equally probable (the best case), then given **only** $5.1 \times 10^9$ different messages there is a probability of more than $0.5$ that two of them have the same digest. We say that we have a probability of more than $1/2$ of *generating a collision*.

## The attack in practice

Alice wants to get Bob to sell his house for five pounds. She needs to get his digital signature on a contract to this effect. She produces a fair contract in which he agrees to sell his house for, say, three hundred and fifty thousand pounds, and an unfair contract in which he agrees to sell for five pounds.

She now pads the two contracts by adding spaces, punctuation and florid legal phrases in lots of different ways producing an enormous number of fair and of unfair contracts. Eventually she hopes a version (contract A) of the fair contract will have the same hash as some version (contract B) of the unfair contract. This is slightly different from the true Birthday paradox and the calculation is more delicate, but essentially the same idea works.



In practice this does not mean that Alice can force Bob to sell the house for five pounds. If she tries to do so Bob can produce the fair contract, and show that his signature fits it too. The real importance of the attack is that once it has been exhibited the signature system is no longer regarded as fully trustworthy.

Modern digital signature keys are chosen to be so long that even this attack is infeasible. Nonetheless the calculations of feasibility rely on the assumption that hash values are evenly distributed. This assumption may be false among the space of messages to be hashed, and at the CRYPT2004 conference collisions were announced in the widely used SHA-0, MD4, MD5, HAVAL-128, and RIPEMD hash algorithms.

## Weakening the discrete logarithm

The same paradox weakens discrete log encryption. Recall that the security of the system relies on the fact that it is hard to find a value $n$ such that $y = x^n \bmod (p)$. Trying a brute force attack we would expect to need to test around $p/2$ values of $n$ to have a $50\%$ chance of finding the correct value.

Instead we can compute values of $x^r$ and $yx^{-s}$. Because of the birthday paradox it takes only about $1.2\sqrt{p}$ steps on average to find a pair $x^r = yx^{-s}$ and so $y = x^{r+s}$ solving the discrete log problem.

Note that computationally we have traded time for storage. The pro-

gramme we write to implement such an attack will not take as long as a brute force attack to run, however it will need a lot of storage (Hard disc or RAM) in which to store the values of $x^r$ and $yx^{-s}$ for comparison. This trade off, of time against space, is a common one in algorithm design.

# Chapter 9

# More cryptography

In this section, we shall consider some more classic ciphers, sometimes referred to as 'pen-and-paper' ciphers. We have seen two such ciphers in the previous section, namely the Caesar shift cipher and the Affine shift cipher.

We are particularly interested in the cryptanalysis of these ciphers and how we might be able to crack an enciphered message with only minimal information relating to the key.

**Keyword Substitution Ciphers**

We have already seen a number of ciphers that are reasonably easy to attack. First of all the Caesar shift cipher ($x \mapsto x + s \bmod(26)$) and then the Affine shift cipher ($x \mapsto sx + t \bmod(26)$). Both of these ciphers are relatively easy to attack using frequency analysis. For the Caesar cipher, knowledge of the enciphered value of 1 character allows us to decipher the entire text, whilst for the Affine cipher, we only need to know the enciphered value of 2 characters.

The Caesar shift and Affine shift ciphers are examples of *simple substitution ciphers*, where each letter of our alphabet is replaced, or substituted, by a unique other letter.

| plaintext | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ciphertext | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T |

The drawback with the Caesar or Affine shifts, is that once we know one (or two in the case of Affine) substituted value, we can calculate them all. Ideally we would be better to have a completely random set of substituted values

| plaintext  | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R |
|------------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ciphertext | X | P | R | A | Z | Q | K | N | S | D | T | H | E | F | W | O | I | Y |

The problem with this is that for the sender and receiver to remember the randomised allocation, will be very difficult. The solution is to use a *keyword*. Suppose we choose the keyword *NUMBER*, we construct our allocation as follows

| plaintext  | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q |
|------------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ciphertext | N | U | M | B | E | R | S | T | V | W | X | Y | Z | A | C | D | F |

List the distinct letters of the keyword, followed by all the rest of the letters, starting from the first unused letter after the final letter of the keyword. To reproduce the entire key, we would only need to remember the keyword. In order to decipher this type of cipher we would need knowledge of all (or at least a large number of) substitutions.

For example, if we encipher the plaintext

*The best kinds of people are warm and kind. They are always there and they never mind. The best kinds of people smile and embrace. They support you with strength and grace.*

with the keyword *NUMBER* we arrive at the ciphertext

*ITEUE HIXVA BHCRD ECDYE NGELN GZNAB XVABI TEPNG ENYLN PHITE GENAB ITEPA EKEGZ VABIT EUEHI XVABH CRDEC DYEHZ VYENA BEZUG NMEIT EPHJD DCGIP CJLVI THIGE ASITN ABSGN ME*

which has frequency table



and although we could guess that 'E' has been enciphered as an 'E', it isn't obvious what the other parts of the key are.

We will briefly look at attacking a substitution cipher later.

**Vigenère Cipher**

Successful use of frequency analysis with the Caesar and Affine shifts ciphers relies on the fact that each occurrence of each character in the plaintext, uses the same function to encipher it. This is an inherent weakness of the cipher and one that we would be best to avoid. It would be more useful from a cryptographic point of view, if we could encipher characters in a different way each time we met them in the plaintext. What if we could

develop a technique to encipher some plaintext, where we use a different shift for each character? For example, if we used a Caesar shift with shift value of 5 on the plaintext

*The best kinds of people are warm and kind. They are always there and they never mind. The best kinds of people smile and embrace. They support you with strength and grace.*

we arrive at the ciphertext

*YMJGJ XYPNS IXTKU JTUQJ FWJBF WRFSI PNSIY MJDFW JFQBF DXYMJ WJFSI YMJDS JAJWR NSIYM JGJXY PNSIX TKUJT UQJXR NQJFS IJRGW FHJYM JDXZU UTWYD TZBNY MXYWJ SLYMF SILWF HJ*

which has frequency table



Frequency Analysis

We could reasonably guess that the shift used was 5 as the most frequent character is now J.

Suppose instead we use the same plaintext but use the sequence of shift values

19 7 4 1 4 18 19 10 8 13 3 18 14 5 15 4 14 15 11 4 0 17 4 22 0 17 12 0 13 3 10 8 13 3 19 7 4 24 0 17 4 0 11 22 0 24 18 19 7 4 17 4 0 13 3 19 7 4 24 13 4 21 4 17 12 8 13 3 19 7 4 1 4 18 19 10 8 13 3 18 14 5 15 4 14 15 11 4 18 12 8 11 4 0 13 3 4 12 1 17 0 2 4 19 7 4 24 18 20 15 15 14 17 19 24 14 20 22 8 19 7 18 19 17 4 13 6 19 7 0 13 3 6 17 0 2 4

(that is we shift the first character by 19, the second by 7 etc)

we arrive at the ciphertext

*MOICI KMUQA GKCKE ICEWI AIISA IYAAG UQAGM OIWAI IAWSA WKMOI IIAAG MOIWA IQIIY QAGMO ICIKM UQAGK CKEIC EWIKY QWIAA GIYCI AEIMO IWKOE ECIMW COSQM OKMII AMMOA AGMIA EI*

which has frequency table

This is very different and it is not at all clear how we can use frequency analysis to decipher the ciphertext.

However, there is an obvious problem! How on earth do we (and the intended recipient) remember the sequence of shift values we need to use? The solution is to use a *keyword* or a *keyphrase*.

As with the keyword substitution ciphers, we pick a memorable keyword or a keyphrase, and we use the numerical values of the characters (A=0, B=1, etc) as the shift value for the corresponding letter.

For the example above, the keyphrase was

*THE BEST KINDS OF PEOPLE ARE WARM AND KIND. THEY ARE ALWAYS THERE AND THEY NEVER MINDTHE BEST KINDS OF PEOPLE SMILE AND EMBRACETHEY SUPPORT YOU WITH STRENGTH AND GRACE*

(this is not a very sensible choice of key, for obvious reasons!)  and the corresponding numerical values were

19 7 4 1 4 18 19 10 8 13 3 18 14 5 15 4 14 15 11 4 0 17 4 22 0 17 12 0 13 3 10 8
13 3 19 7 4 24 0 17 4 0 11 22 0 24 18 19 7 4 17 4 0 13 3 19 7 4 24 13 4 21 4 17
12 8 13 3 19 7 4 1 4 18 19 10 8 13 3 18 14 5 15 4 14 15 11 4 18 12 8 11 4 0 13
3 4 12 1 17 0 2 4 19 7 4 24 18 20 15 15 14 17 19 24 14 20 22 8 19 7 18 19 17 4
13 6 19 7 0 13 3 6 17 0 2 4

Of course in general we would normally use a much smaller keyphrase that is shorter then the length of the original plaintext!

Suppose we had used the keyword Caesar on the text

*According to Suetonius, Caesar simply replaced each letter in a message with the letter that is three places further down the alphabet. Cryptographers often think in terms of the plaintext alphabet as being the alphabet used to write the original message, and the ciphertext alphabet as being the letters that are substituted in place of the plain letters. When the plaintext alphabet is placed above the ciphertext alphabet, as shown below, it is clear to see that the ciphertext alphabet has been shifted by three places. Hence this form of substitution is often called the Caesar Shift Cipher. A cipher is the name given to any form of cryptographic substitution, in which each letter is replaced by another letter or symbol.*

The shift values we would then use are

2 0 4 18 0 17

and we simply repeat this pattern as often as necessary.

2 0 4 18 0 17 2 0 4 18 0 17 2 0 4 18 0 17 2 0 4 18 0 17 ……. We arrive at the ciphertext

*CCGGR UKNKL OJWEX GNZWS GSEJC RWAMG NYVWP CCCIV EREHP WTKGR MFADG SWSGV YIXZT YGLIL TVTTL STZUT LJEVR LEUEJ HUVLH VTDSO NKJEE DPYCB ILCIA PXGGI CPLWR JQFXW NKJIR CIEVE VESFH TLWPC CIRLE OVAPH HRDEX SSSGI RYTYG APHHR DEXMS VFTSO RZVEX ZEFTI KANRN MIKSR IEEFD KJEGA PYGRX WXKCL TZASG TEKBV KN-KLH VNEXL EIUTL STRTE WMBJV IXMTV FIRHL REESX TYGPP SIENE XLEIU WLWNK JETDA ZPTIP TRNPL SBVVI WHLRE EHSBF XEXZE TK-PLW RKGXX SLGJA FWTRU SLGWE DEPGW ZVIWU LVCRX GSVGT LSTKJ EGAPY GRXWX KCLTZ ASGTL SSSGE RKHZH TIVBP VHVWE GNAGW SYGNG WTYKS JGRDQ FWMBJ VIXMT ZQNMK OWVER UACNE HLHVE AIKAI UHMXT TKPLW RREIT ZEIKS XZEEC MIYIM GNXGA EAFSJ MFHCV QPKQG VSPYK CWMBJ VIXMT ZQNMF WYKCL WATJL ILTVT IWJEG NAGWD SAARG TYGRP WTKGR SJSPO BSDTY KSXWX KKSXS KVPFV GMYVT TOWNU IQGNJ KNKZN VVTLW BCCCO UHROB IJCRG SEJHK OL*

with frequency table



Frequency Analysis

To assist the process we can make use of the following table:

```
  A B C D E F G H I J K L M N O P Q R S T U V W X Y Z
A A B C D E F G H I J K L M N O P Q R S T U V W X Y Z
B B C D E F G H I J K L M N O P Q R S T U V W X Y Z A
C C D E F G H I J K L M N O P Q R S T U V W X Y Z A B
D D E F G H I J K L M N O P Q R S T U V W X Y Z A B C
E E F G H I J K L M N O P Q R S T U V W X Y Z A B C D
F F G H I J K L M N O P Q R S T U V W X Y Z A B C D E
G G H I J K L M N O P Q R S T U V W X Y Z A B C D E F
H H I J K L M N O P Q R S T U V W X Y Z A B C D E F G
I I J K L M N O P Q R S T U V W X Y Z A B C D E F G H
J J K L M N O P Q R S T U V W X Y Z A B C D E F G H I
K K L M N O P Q R S T U V W X Y Z A B C D E F G H I J
L L M N O P Q R S T U V W X Y Z A B C D E F G H I J K
M M N O P Q R S T U V W X Y Z A B C D E F G H I J K L
N N O P Q R S T U V W X Y Z A B C D E F G H I J K L M
O O P Q R S T U V W X Y Z A B C D E F G H I J K L M N
P P Q R S T U V W X Y Z A B C D E F G H I J K L M N O
Q Q R S T U V W X Y Z A B C D E F G H I J K L M N O P
R R S T U V W X Y Z A B C D E F G H I J K L M N O P Q
S S T U V W X Y Z A B C D E F G H I J K L M N O P Q R
T T U V W X Y Z A B C D E F G H I J K L M N O P Q R S
U U V W X Y Z A B C D E F G H I J K L M N O P Q R S T
V V W X Y Z A B C D E F G H I J K L M N O P Q R S T U
W W X Y Z A B C D E F G H I J K L M N O P Q R S T U V
X X Y Z A B C D E F G H I J K L M N O P Q R S T U V W
Y Y Z A B C D E F G H I J K L M N O P Q R S T U V W X
Z Z A B C D E F G H I J K L M N O P Q R S T U V W X Y
```

This is known as a *tabula recta*. The letters in the first row correspond to the letters in the plaintext, while the letters in the first column correspond to the letters in the keyphrase.

For each letter in the keyphrase, lookup the corresponding row in the tabula recta and then look along for the corresponding letter in the plaintext

| PLAINTEXT  | A | C | C | O | R | D | I | N | G | T | O | S | .... |
|------------|---|---|---|---|---|---|---|---|---|---|---|---|------|
| KEYPHRASE  | C | A | E | S | A | R | C | A | E | S | A | R | ...  |
| CIPHERTEXT | C | C | G | G | R | U | K | N | K | L | O | J | ...  |

In practice you may find it easier to split your plaintext into columns - 1 column for each letter of the keyword.

| C | A | E | S | A | R |
|---|---|---|---|---|---|
| A | C | C | O | R | D |
| I | N | G | T | O | S |
| U | E | T | O | N | I |
| U | S | C | A | E | S |
| A | R | S | I | M | P |
| L | Y | R | E | P | L |
| A | C | E | D | E | A |
| C | H | L | E | T | T |
| E | R | I | N | A | M |
| E | S | S | A | G | E |
| W | I | T | H | T | H |
| E | L | E | T | T | E |
| R | T | H | A | T | I |
| S | T | H | R | E | E |
| P | L | A | C | E | S |
| F | U | R | T | H | E |
| R | D | O | W | N | T |
| H | E | A | L | P | H |
| A | B | E | T |   |   |

Each column is then enciphered using the same Caesar shift, the value of which is given by the encoded value of the corresponding character in the keyphrase.

This type of cipher is known as a Vigenère Cipher, named after Blaise de Vigenère, who published a description of a similar cipher in 1586. The actual cipher we have described was really invented by Giovan Battista Bellaso, who built on the work of Leon Battista Alberti and Johannes Trithemius. It is an example of what is referred to as a polyalphabetic cipher where multiple substitution alphabets are used.

### Keyword Transposition Ciphers

With a *transposition cipher*, rather than replace characters by other characters according to some rule, we simply transpose the order of the characters in the plaintext. Hence the ciphertext becomes simply a permutation of the plaintext. A simple example of this is the scytale that we saw at the start of Chapter 8. A more sophisticated example uses a keyword in the following way:

1. Choose a keyword (e.g. TRANSPOSITION) and strip off any repeated letters (e.g. TRANSPOI)
2. Form a grid with the reduced keyword in the first row and then the rest of the plaintext in the subsequent rows. E.g.

| T | R | A | N | S | P | O | I |
|---|---|---|---|---|---|---|---|
| A | C | C | O | R | D | I | N |
| G | T | O | S | U | E | T | O |
| N | I | U | S | C | A | E | S |
| A | R | S | I | M | P | L | Y |
| R | E | P | L | A | C | E | D |
| E | A | C | H | L | E | T | T |
| E | R | I | N | A | M | E | S |
| S | A | G | E | W | I | T | H |
| T | H | E | L | E | T | T | E |
| R | T | H | A | T | I | S | T |
| H | R | E | E | P | L | A | C |
| E | S | F | U | R | T | H | E |
| R | D | O | W | N | T | H | E |
| A | L | P | H | A | B | E | T |

3. If the plaintext does not fit into the grid exactly, then we can, optionally, use 'padding characters' to fill in any gaps at the end. 4. Now read down each column in turn, starting with the column whose first letter (i.e. the letter from the keyword) comes first in the alphabet and then proceeding alphabetically from then on.

Using the above example, we get the ciphertext

*COUSP CIGEH EFOPN OSYDT SHETC EETOS SILHN ELAEU WHITE LETET TSAHH EDEAP CEMIT ILTTB CTIRE ARAHT RSDLR UCMAL AWETP RNAAG NAREE STRHE RA*

## 9.1   Deciphering the Vigenère cipher

Now that we know how a Vigenère cipher works, how do we attack the cipher?

At first sight, frequency analysis doesn't seem to be any use to us. However, notice that if we knew the length of the keyphrase, we could reconstruct the columns above (using the ciphertext) and use frequency analysis to determine the key letter used for each column. So all we have to do is learn the length of the keyphrase.

There are two common techniques for this.

**Kasiski Test**

The idea, put forward independently by Babbage and Kasiski, is that in a large piece of plaintext, it may happen that some common words, or groups of letters, may appear a number of times starting at the same column in the

columnar description above. If this is the case, the distance between these words will be a multiple of the length of the keyphrase. Some of these occurrences will of course be purely random but some may not. Hence we look for some common groups of letters in the ciphertext and determine the distance between them. This will give us clues to the length of the keyphrase.

Consider the previous example of a Vigenère cipher

*CCGGR UKNKL OJWEX GNZWS GSEJC RWAMG NYVWP CCCIV EREHP WTKGR MFADG SWSGV YIXZT YGLIL TVTTL STZUT LJEVR LEUEJ HUVLH VTDSO NKJEE DPYCB ILCIA PXGGI CPLWR JQFXW NKJIR CIEVE VESFH TLWPC CIRLE OVAPH HRDEX SSSGI RYTYG APHHR DEXMS VFTSO RZVEX ZEFTI KANRN MIKSR IEEFD KJEGA PYGRX WXKCL TZASG TEKBV KN-KLH VNEXL EIUTL STRTE WMBJV IXMTV FIRHL REESX TYGPP SIENE XLEIU WLWNK JETDA ZPTIP TRNPL SBVVI WHLRE EHSBF XEXZE TK-PLW RKGXX SLGJA FWTRU SLGWE DEPGW ZVIWU LVCRX GSVGT LSTKJ EGAPY GRXWX KCLTZ ASGTL SSSGE RKHZH TIVBP VHVWE GNAGW SYGNG WTYKS JGRDQ FWMBJ VIXMT ZQNMK OWVER UACNE HLHVE AIKAI UHMXT TKPLW RREIT ZEIKS XZEEC MIYIM GNXGA EAFSJ MFHCV QPKQG VSPYK CWMBJ VIXMT ZQNMF WYKCL WATJL ILTVT IWJEG NAGWD SAARG TYGRP WTKGR SJSPO BSDTY KSXWX KKSXS KVPFV GMYVT TOWNU IQGNJ KNKZN VVTLW BCCCO UHROB IJCRG SEJHK OL*

We have the following collection of 'words' of length 3 and their relative distances apart.

| Substring | Positions | Intervals |
|-----------|-----------|-----------|
| KNK | 7,241,631 | 234,390 |
| NKL | 8,242 | 234 |
| WSG | 19,57 | 38 |
| GSE | 21,655 | 634 |
| WPC | 34,148 | 114 |
| PCC | 35,149 | 114 |
| CCC | 36,642 | 606 |
| CCI | 37,150 | 113 |
| VER | 40,463 | 423 |
| PWT | 45,585 | 540 |
| WTK | 46,586 | 540 |
| TKG | 47,587 | 540 |
| KGR | 48,588 | 540 |

Now take these intervals and calculate their greatest common divisors. The keylength is likely to be a divisor of one of these numbers.

| | 234 | 390 | 38 | 634 | 114 | 606 | 423 | 540 |
|---|---|---|---|---|---|---|---|---|
| 234 | | 78 | 2 | 2 | 6 | 6 | 9 | 18 |
| 390 | | | 2 | 2 | 6 | 6 | 3 | 30 |
| 38 | | | | 2 | 38 | 2 | 1 | 2 |
| 634 | | | | | 2 | 2 | 1 | 2 |
| 114 | | | | | | 6 | 3 | 6 |
| 606 | | | | | | | 3 | 6 |
| 113 | | | | | | | 1 | 1 |
| 423 | | | | | | | | 9 |

The values of 2 and 6 look likely candidates for the keylength.

Once we have determined the keylength, then we can use frequency analysis on each column to try to determine which shift was used. We shall look into this in more details later.

The **index of coincidence** is a way to measure how close to being English a particular piece of text is. It relies on frequency analysis and a longer explanation of it can be found later, (see 9.2). The important property for us is that it is invariant under a simple substitution cipher such as a Caesar shift. Essentially we try different values for the keyword length and measure how close to being English the corresponding text in each column is. The highest value is likely to be (a multiple of) the column length.

For example, for the ciphertext

*CCGGR UKNKL OJWEX GNZWS GSEJC RWAMG NYVWP CCCIV EREHP WTKGR MFADG SWSGV YIXZT YGLIL TVTTL STZUT LJEVR LEUEJ HUVLH VTDSO NKJEE DPYCB ILCIA PXGGI CPLWR JQFXW NKJIR CIEVE VESFH TLWPC CIRLE OVAPH HRDEX SSSGI RYTYG APHHR DEXMS VFTSO RZVEX ZEFTI KANRN MIKSR IEEFD KJEGA PYGRX WXKCL TZASG TEKBV KN-KLH VNEXL EIUTL STRTE WMBJV IXMTV FIRHL REESX TYGPP SIENE XLEIU WLWNK JETDA ZPTIP TRNPL SBVVI WHLRE EHSBF XEXZE TK-PLW RKGXX SLGJA FWTRU SLGWE DEPGW ZVIWU LVCRX GSVGT LSTKJ EGAPY GRXWX KCLTZ ASGTL SSSGE RKHZH TIVBP VHVWE GNAGW SYGNG WTYKS JGRDQ FWMBJ VIXMT ZQNMK OWVER UACNE HLHVE AIKAI UHMXT TKPLW RREIT ZEIKS XZEEC MIYIM GNXGA EAFSJ MFHCV QPKQG VSPYK CWMBJ VIXMT ZQNMF WYKCL WATJL ILTVT IWJEG NAGWD SAARG TYGRP WTKGR SJSPO BSDTY KSXWX KKSXS KVPFV GMYVT TOWNU IQGNJ KNKZN VVTLW BCCCO UHROB IJCRG SEJHK OL*

we can compute the Index for different values of keyphrase length and get

Index of Coincidence

This has 'spikes' at 6, 12 and 18 suggesting that the keyword might have length 6. Frequency analysis on each of the 6 columns would then yield the plaintext.

## 9.2 Index of Coincidence

The index of coincidence is a measure of how close a frequency distribution is to the uniform distribution. If you choose two letters at random from a random piece of text, the probability that they are the same is about 0.0385, whereas if you choose two letters at random from a piece of English text, the probability that they are the same if about 0.0667. This probability does not change if the text is enciphered with a substitution cipher. We can therefore exploit this phenomenon to decide if a given piece of ciphertext has been enciphered by a substitution cipher such as a Caesar shift.

Given a piece of text of length $N$ in which the letters 'A' to 'Z' appear with frequency $f_i$ ($i = A$ to $Z$) then the index of coincidence of the text is defined as

$$IC = \frac{26 \sum_{i=A}^{Z} f_i(f_i - 1)}{N(N - 1)}.$$

Note that some definitions omit the value 26 from the above definition.

If the text is essentially 'random' then the index of coincidence will be close to $1$. If it is text written in English (and possibly enciphered using a substitution cipher) then it will be closer to $1.734 (= 26 \times 0.0667)$.

We can exploit this when analysing a piece of ciphertext that we suspect has been encoded using a vigenère cipher. If we have a guess that the keylength is $n$ then we can write the ciphertext using $n$ columns and extract all $n$ columns from the ciphertext. Given that each of these has been enciphered using a Caesar shift (with a different value for the shift), we can calculate the index of coincidence for each column and if these are (or better, their average is) close to $1.734$, then there is a good chance that $n$ really is the key length. If not, then we can try another value of $n$.

For example the following text has been enciphered using a Vigenère cipher with the keyword *CIPHERS*

*KVRYC GLQOG HTYQC KXWLV JQZRF TYWTQ HHRRD IWGPX YEHWG WIIXQ ZBPRX WPKGF TKAQV DYHVU TGEAM FFKVP ZIIAG ADMAV DNLTM MEWFA ILTJL JIIJE ETGND SPFOG LPZEG JQKTK YIWCV PSXVJ PIIPZ VDGAH JSDEQ VILVD AUMCJ MGZGZ BLRKL QMCJM GZGZD YIEUQ LTPWK GEWCC IILKV UVVDS VQDUM ELQKX WLVJQ ZRVHV LCSTU JIGOP IATJW PEXRM GWFQP VVXOK SXJMG ZGZ*

If we calculate the index of coincidence based on the text having a key length of either 1 or 2 or 3 or 4 or ... or 15, then we get the following chart



Index of Coincidence

There is clearly a large 'spike' at the value 7 (and 14) which indicates that the length of the keyword is probably 7. Frequency analysis on the 7 columns can then be used to try to decipher the text.

(in case you are wondering, the plaintext is

*In cryptography, a cipher (or cypher) is an algorithm for performing encryption or decryption in a series of well-defined steps that can be followed as a procedure. An alternative, less common term is encipherment. To encipher or encode is to convert information into cipher or code. (taken from https://en.wikipedia.org/wiki/Cipher)*

with all the non-alphabetic characters stripped out.)

## 9.3   Chi-squared test

As an alternative to the index of coincidence, we can use a technique that is dependant on a statictical test, called the *Chi-squared test*. This is a measure of how similar a distribution is to the expected values for that type of distribution. In cryptography we define the *Chi-squared statistic* as

$$\chi^2(C, E) = \sum_{i=A}^{i=Z} \frac{(C_i - E_i)^2}{E_i}$$

where $C_i$ is the actual count of the letter represented by $i$ in our text, while $E_i$ is the expected count (taking account of the expected frequency of let-

ters from whichever natural language the text is taken from - in our case, normally English) (in the table below, we 'normalize' the values by dividing by $n$, where $n$ is the length of the text). In essence, the smaller the Chi-squared statistic, the more likely it is to have been written in English.

As an example, the text

*FTQNQ EFWUZ PEARB QABXQ MDQIM DYMZP WUZPF TQKMD QMXIM KEFTQ DQMZP FTQKZ QHQDY UZPFT QNQEF WUZPE ARBQA BXQEY UXQMZ PQYND MOQFT QKEGB BADFK AGIUF TEFDQ ZSFTM ZPSDM OQ*

is the output of a caesar shift. By considering all possible shift values, we can see that the lowest chi-squared statistic occurs for a shift of 12, suggesting that this is the shift used.



I leave it as an exercise for you to decipher the text.

## 9.4 Frequency Analysis



If you encipher a plaintext message using a substitution cipher then the letters in the ciphertext will occur with the corresponding frequency of their plaintext values. This information enables us to attack the ciphertext using frequency analysis. One counts the number of occurrences of each character in the ciphertext and compares it with an expected frequency for the standard English alphabet.

For example if we consider the following text which has been enciphered using a Caesar shift cipher

*FHHTW INSLY TXZJY TSNZX HFJXF WXNRU QDWJU QFHJI JFHMQ JYYJW NSFRJ XXFLJ BNYMY MJQJY YJWYM FYNXY MWJJU QFHJX KZWYM JWITB SYMJF QUMFG JYHWD UYTLW FUMJW XTKYJ SYMNS PNSYJ WRXTK YMJUQ FNSYJ CYFQU MFGJY FXGJN SLYMJ FQUMF GJYZX JIYTB WNYJY MJTWN LNSFQ RJXXF LJFSI YMJHN UMJWY JCYFQ UMFGJ YFXGJ NSLYM JQJYY JWXYM FYFWJ XZGXY NYZYJ INSUQ FHJTK YMJUQ FNSQJ YYJWX BMJSY MJUQF NSYJC YFQUM FGJYN XUQFH JIFGT AJYMJ HNUMJ WYJCY FQUMF GJYFX XMTBS GJQTB NYNXH QJFWY TXJJY MFYYM JHNUM JWYJC YFQUM FGJYM FXGJJ SXMNK YJIGD YMWJJ UQFHJ XMJSH JYMNX KTWRT KXZGX YNYZY NTSNX TKYJS HFQQJ IYMJH FJXFW XMNKY HNUMJ WFHNU MJWNX YMJSF RJLNA JSYTF SDKTW RTKHW DUYTL WFUMN HXZGX YNYZY NTSNS BMNHM JFHMQ JYYJW NXWJU QFHJI GDFST YMJWQ JYYJW TWXDR GTQYM NXYJC YNXYF PJSKW TRMYY UBBBX NRTSX NSLMS JYYMJ GQFHP HMFRG JWHFJ XFWMY RQ*

The individual characters appear with the following frequencies



We can compare this to the chart above and guess that the letter E, the most popular letter in English text, is encoded as J which would give a shift value of 5. This is indeed correct.

**ASINTOER Frequency**

Indeed, if we take any substring in the ciphertext, say for example every second character, then the letter frequencies for this substring should be approximately the same as for the original string. For example, taking every 4th character in the previous ciphertext, starting at character 4, we get

*TSXTX XNDQI MYSXJ MQJFY JHZJB JMYUW JKYPJ TJNCU JGLFF ZYNMN FXJYN WYMYJ YJWFJ XZNFK USYBY QYFFN FFJHJ CUJXS TNJTY YNWYM YGXYD JFMJX RZNNX JQYFW KUFMX SLSSW HYFHX ZSMJQ JWFGT WYWGM JXJTY BTSJJ HFWXY*

with frequency table

The letter J is still the most frequent character by a long way and the two charts are very similar in nature. We could reasonably guess from this that the shift was 5.

However, this does not always hold. For example, taking every 5th character in the string above starting at character 3, we get

*HSZNJ NWHHY FFYQW NJHWI MMHTM KMSXJ SFGGY UYYYT SXFJJ YFXLJ XFGZS JJSJJ UYQJQ FYUJU YTQNF JYNYQ JGMIW FJMWZ YSYQM XNUHW JLYKK YUZYS NHYWH FJYXQ YXSMB TLYFF HW*

with frequency table



The frequency chart is a little bit different to the previous one, and the most frequent character now is Y, which would give a shift value of 20.

In this case, instead of just considering the single most common letter in English, we can consider multiple letters at the same time. For example, 8 of the most common letters in the English language are E, T, A, O, I, N, S and R. In the following chart, we have considered each of the 26 possible shift values (0 - 25) and calculated the sum of the frequencies of the enciphered values of the five letters E, T, A, O, I, N, S and R (we use the mnemonic A SIN TO ER to remember these).



It is then clear that the highest aggregate frequency occurs for a shift of 5, confirming the shift value for the original text. This technique can be

particularly useful when deciphering the Vigenère cipher. Note however that it might not always be accurate! You have to be prepared to use some trial and error when analysing a ciphertext.

In Appendix H we have set up a form to encrypt and decrypt Vigenère ciphers and added a version of a cracking algorithm for you to play with. In our algorithm we use the index of coincidence to 'guess' the length of the keyword (calculate the index of coincidence for lots of possible keylengths and then choose the 'best'), and then use ASINTOER frequency analysis to determine the shifted value for each column, and so recover the keyword.

## 9.5   Cracking Keyword Substitution Cipher

Attacking a keyword substitution cipher is not always easy. Although we can benefit from frequency analysis, it may only be good enough to give us a few elements from the key, and a long and complicated key may well cause serious problems.

There is however another technique, which can also be used in solving other ciphers, that we can employ, and is known as *Hill Climbing*.



The idea is as follows: for each potential key we assign a score as to how suitable that key is to be the correct key (we can use any suitable technique we like, such as some other type of frequency analysis). If small changes in the key result in small changes in the score, then we can then think of the results as providing us with a graph with the higher values representing more likely keys. We then 'guess' an initial key and work out its score. We then change the key slightly and recalculate the score. If this new score is 'worse' than the original, then we continue to make small changes to the key until we either get a 'better' key or run out of possible changes to the key. If we get a better key, then we adopt that one as our current best bet and start the process over again, otherwise we stop the algorithm and our present key is then our best bet.

Hopefully, this will enable us to find the local maximum of the graph of scores, and with some luck, this local maximum may provide us with the correct key.

In Appendix I we have set up a form to encrypt and decrypt keyword substitution ciphers and added a version of the hill-climbing algorithm for you to play with.In our algorithm we use frequency analysis based on the frequency of *trigrams*, groups of three consecutive letters, but analysis using quadgrams (groups of 4 letters) is usually preferred as being more accurate.

# Bibliography

Jones, G. A. and Jones, J. M. (2006). *Elementary number theory*. Springer (India), New Delhi.

Liebek, M. (2015). *A Concide Introduction to Pure Mathematics*. Chapman & Hall.

Rosen, K. H. (2010). *Elementary Number Theory*. Pearson.

Solow, D. (2013). *How to read and do proofs: An introduction to mathematical thought processes*. John Wiley & Sons.

Velleman, D. (2006). *How to prove it*. Cambridge University Press.

# Index