

# **Elementary Number Theory**

**ETH Zürich**

Dr. Markus Schwagenscheidt

Spring Term 2023 - Last update: October 3, 2023



# Contents

<b>1</b>	<b>Primes and Divisibility</b>	<b>5</b>
1.1	Prime numbers and the Fundamental Theorem of Arithmetic . . . . .	5
1.2	The greatest common divisor and the Euclidean algorithm . . . . .	7
1.3	Distribution of primes . . . . .	10
<b>2</b>	<b>Number-theoretic functions</b>	<b>19</b>
2.1	Multiplicative number-theoretic functions . . . . .	19
2.2	Dirichlet convolution and Möbius inversion . . . . .	21
2.3	Perfect and amicable numbers . . . . .	25
<b>3</b>	<b>Modular arithmetic</b>	<b>29</b>
3.1	Congruences . . . . .	29
3.2	Fermat's Little Theorem . . . . .	31
3.3	The theorems of Lagrange, Wilson, and Wolstenholme . . . . .	33
3.4	Quadratic residues . . . . .	36
3.5	Applications to cryptography . . . . .	41
<b>4</b>	<b>Quadratic forms</b>	<b>49</b>
4.1	Sums of squares . . . . .	49
4.2	Binary quadratic forms . . . . .	51
4.3	Reduction theory of binary quadratic forms . . . . .	54
4.4	Gauss composition . . . . .	57
4.5	Pell's equation and continued fractions . . . . .	60
4.6	Congruent numbers . . . . .	69
<b>5</b>	<b>Partitions</b>	<b>77</b>
5.1	The partition function . . . . .	77
5.2	Generating functions and Euler's Pentagonal Theorem . . . . .	79



# 1 Primes and Divisibility

The set of *natural numbers* is

$$\mathbb{N} = \{1, 2, 3, 4, \dots\},$$

and the set of *integers* is

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

We also put

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

We will use the following properties of the integers:

1. *Associativity*:  $a + (b + c) = (a + b) + c$  and  $a(bc) = (ab)c$ .
2. *Commutativity*:  $a + b = b + a$  and  $ab = ba$ .
3. *Distributivity*:  $a(b + c) = ab + ac$ .
4. *Identity*:  $a + 0 = 0 + a = a$  and  $a \cdot 1 = 1 \cdot a = a$ .
5. *Inverse*:  $a + (-a) = (-a) + a = 0$ .
6. *Transitivity*:  $a > b$  and  $b > c$  implies  $a > c$ .
7. *Trichotomy*: Either  $a > b$ ,  $a < b$  or  $a = b$ .
8. *Cancellation law*: If  $ac = bc$  and  $c \neq 0$ , then  $a = b$ .

Moreover, we will use that the natural numbers are *well-ordered*: every non-empty subset of  $\mathbb{N}$  has a smallest element.

## 1.1 Prime numbers and the Fundamental Theorem of Arithmetic

**Definition 1.1.1.** For  $a, b \in \mathbb{Z}$  we say that  $a$  *divides*  $b$ , written  $a \mid b$ , if  $ac = b$  for some  $c \in \mathbb{Z}$ . In this case, we say that  $a$  is a *divisor* of  $b$ . If  $a$  is not a divisor of  $b$ , then we say that  $a$  *does not divide*  $b$ , and write  $a \nmid b$ .

For example,  $2 \mid 4$  and  $-5 \mid 15$  and  $-3 \mid -18$ . We have the following fundamental rules for divisibility, whose proof we leave as an exercise.

**Lemma 1.1.2.**

1.  $a \mid 0$ ,  $1 \mid a$ , and  $a \mid a$  for all  $a \in \mathbb{Z}$ ,
2.  $a \mid b$  implies  $|a| \leq |b|$  if  $b \neq 0$ ,
3.  $a \mid b$  and  $b \mid a$  imply  $a = \pm b$ ,

## 1 Primes and Divisibility

4.  $a \mid b$  and  $b \mid c$  imply  $a \mid c$ ,

5.  $a \mid b_1$  and  $a \mid b_2$  imply  $a \mid (mb_1 + nb_2)$  for all  $m, n \in \mathbb{Z}$ ,

6.  $a_1 \mid b_1$  and  $a_2 \mid b_2$  imply  $a_1a_2 \mid b_1b_2$ .

**Definition 1.1.3.** An integer  $n > 1$  is *prime* if the only positive divisors of  $n$  are 1 and  $n$ . We call  $n > 1$  *composite* if  $n$  is not prime.

Note that 1 is neither prime nor composite by our definition. The first few primes are

$$2, 3, 5, 7, 11, 13, 17, 19, 23, 29, 31, 37, 41, 43, 47, 53, 59, 61, 67, 71, 73, 79, \dots$$

The prime numbers form the building blocks of all natural numbers in the following sense.

**Theorem 1.1.4** (Fundamental Theorem of Arithmetic). *Every natural number  $n > 1$  can be written as a product of primes,*

$$n = p_1 \cdots p_r,$$

*uniquely up to order, with (not necessarily pairwise different) primes  $p_1, \dots, p_r$ . This is called the prime factorization of  $n$ .*

**Remark 1.1.5.** We usually collect the powers of the same prime in the prime factorization of  $n$  and write

$$n = p_1^{\nu_1} \cdots p_r^{\nu_r}$$

with *pairwise different* primes  $p_1, \dots, p_r$  and *multiplicities*  $\nu_1, \dots, \nu_r \in \mathbb{N}$ . For a prime  $p$  and a natural number  $n > 1$ , the multiplicity (or *p-adic valuation*) of  $p$  in  $n$  is also denoted by  $\nu_p(n)$ . We set  $\nu_p(1) = 0$ . Note that  $\nu_p(n) > 0$  if and only if  $p$  divides  $n$ , so we can write the prime factorization of  $n$  as

$$n = \prod_p p^{\nu_p(n)},$$

where the product runs over all primes. The *prime omega functions*  $\omega(n)$  and  $\Omega(n)$  count the number of different prime factors of  $n$  (with and without multiplicities, respectively), that is, for  $n = p_1^{\nu_1} \cdots p_r^{\nu_r}$  they are defined by

$$\omega(n) = r, \quad \Omega(n) = \nu_1 + \cdots + \nu_r.$$

The uniqueness of the prime factorization requires Euclid's Lemma, whose proof we postpone to Lemma 1.2.6 below as it needs some facts about the greatest common divisor.

**Lemma 1.1.6** (Euclid). *Let  $p$  be a prime and  $a, b \in \mathbb{N}$ . If  $p \mid ab$  then  $p \mid a$  or  $p \mid b$ .*

*Proof of Theorem 1.1.4. Existence of the prime factorization.* The proof proceeds by induction over  $n \in \mathbb{N}$  with  $n > 1$ . For  $n = 2$  the claim is clear since 2 is a prime number. If  $n > 2$  is itself a prime number, then we are done, and if  $n$  is composite, then we may write  $n = ab$  with  $1 < a, b < n$ . By induction hypothesis,  $a$  and  $b$  have prime factorizations, so  $n$  has a prime factorization as well.

*Uniqueness of the prime factorization.* Let  $n \in \mathbb{N}$  with  $n > 1$ , and let

$$n = p_1 p_2 \cdots p_r, \quad n = q_1 q_2 \cdots q_s,$$

## 1.2 The greatest common divisor and the Euclidean algorithm

be two prime factorizations of  $n$ , where the  $p_i$  or  $q_j$  need not be pairwise different. By Euclid's Lemma we know that  $p_1 \mid q_1$  or  $p_1 \mid q_2 \cdots q_s$ . If  $p_1 \mid q_1$  then  $p_1 = q_1$  since both are prime, and if  $p_1 \mid q_2 \cdots q_s$  then we can repeat the argument until we find that  $p_1 = q_k$  for some  $k \in \{1, \dots, s\}$ . Hence, by reordering the factors  $q_j$  we can assume that  $p_1 = q_1$ . Repeating this process for the smaller number  $n' = p_2 \cdots p_r = q_2 \cdots q_s$ , we obtain the uniqueness of the prime factorization.  $\square$

Using the existence of the prime factorization one can show that there are infinitely many primes.

**Theorem 1.1.7** (Euclid). *There are infinitely many primes.*

*Proof.* We leave the proof as an exercise, with the following hint: suppose that there are only finitely many distinct primes  $p_1, p_2, \dots, p_r$ , and consider prime factorization of the natural number  $m = p_1 p_2 \cdots p_r + 1$ .  $\square$

It is often useful to have an infinite supply of primes with certain additional properties. A very important result in this direction, whose proof goes beyond the methods of this lecture, is due to Dirichlet.

**Theorem 1.1.8** (Dirichlet's Theorem on Primes in Arithmetic Progressions). *For  $a, b \in \mathbb{N}$  with  $\gcd(a, b) = 1$  there exist infinitely many primes of the form  $am + b$ , with  $m \in \mathbb{N}$ .*

For example, we will show in the exercises that there are infinitely many primes of the form  $4n + 3$  (e.g. 3, 7, 11, 19, 23, ...).

A useful method to list all primes up to a given bound  $N$  is given by the *Sieve of Eratosthenes*, which works as follows: make a list of all integers from 2 to  $N$ . The smallest number, 2, must be a prime. Now every second number (4, 6, 8, ...) is even, hence not prime, so cross them out. The smallest integer greater than 2 not crossed out is 3, so it must be prime. Every third number (6, 9, 12, ...) is divisible by 3, so cross them out. The next number not crossed out is 5, so it must be prime. Cross out every fifth number (10, 15, 20, ...). We continue with this process, and once the largest prime less than  $\sqrt{N}$  is reached, the list will only contain primes.

## 1.2 The greatest common divisor and the Euclidean algorithm

**Definition 1.2.1.** For  $a, b \in \mathbb{Z}$  not both 0 we call the natural number

$$\gcd(a, b) = \max\{d \in \mathbb{N} : d \mid a \text{ and } d \mid b\}$$

the *greatest common divisor* (or *gcd*) of  $a$  and  $b$ . We put  $\gcd(0, 0) = 0$ . We call  $a$  and  $b$  *coprime* if  $\gcd(a, b) = 1$ .

**Remark 1.2.2.** We can also define the gcd of  $r$  integers  $a_1, \dots, a_r$  (not all 0) by

$$\gcd(a_1, \dots, a_r) = \max\{d \in \mathbb{N} : d \mid a_j \text{ for all } j = 1, \dots, r\}.$$

We call  $a_1, \dots, a_r$  *coprime* if  $\gcd(a_1, \dots, a_r) = 1$ , and we call them *pairwise coprime* if  $\gcd(a_i, a_j) = 1$  for all  $i \neq j$ . Since

$$\gcd(a_1, \dots, a_r) = \gcd(\gcd(a_1, \dots, a_{r-1}), a_r),$$

we mostly restrict ourselves to the gcd of two numbers for notational simplicity.

## 1 Primes and Divisibility

In order to compute the gcd, we will use division with remainder.

**Proposition 1.2.3** (Division with remainder). *Let  $a, b \in \mathbb{Z}$  with  $b \neq 0$ . Then there exist unique  $q, r \in \mathbb{Z}$  with  $0 \leq r < |b|$  such that*

$$a = bq + r.$$

*Proof.* For simplicity, we assume that  $b > 0$ , but the proof for  $b < 0$  is analogous. We consider the set

$$M = \{r \in \mathbb{N}_0 : r = a - qb, q \in \mathbb{Z}\}.$$

Since  $b > 0$  we have  $a + |a|b \in M$ , so the set  $M$  is non-empty and hence has a smallest element  $r = a - bq$ . The minimality of  $r$  implies  $0 \leq r < b$ .

Suppose that  $q', r' \in \mathbb{Z}$  satisfy  $0 \leq r' < b$  and  $a = q'b + r'$ . Then  $(q - q')b + (r - r') = 0$ . If  $r \neq r'$ , then  $b \leq |r - r'|$ , which is a contradiction to  $0 \leq r, r' < b$ . This implies  $r = r'$  and then also  $q = q'$ , which shows the uniqueness of  $q$  and  $r$ .  $\square$

It is a useful fact that  $\gcd(a, b)$  can be written as an integral linear combination of  $a$  and  $b$ :

**Lemma 1.2.4** (Bézout's Lemma). *For  $a, b \in \mathbb{Z}$  there exist  $x, y \in \mathbb{Z}$  such that*

$$\gcd(a, b) = ax + by.$$

*Proof.* The statement is clear for  $a = b = 0$ , so we can assume that  $a, b$  are not both 0. We consider the set

$$M = \{ax + by : x, y \in \mathbb{Z}, ax + by > 0\} \subset \mathbb{N},$$

and show that  $\gcd(a, b) = \min M$ . For brevity, we will put  $d = \gcd(a, b)$ . Choosing  $x = a$  and  $b = y$  we see that the set  $M$  is non-empty, and hence has a minimum, which we denote by  $d'$ . We may write  $d' = ax + by$  for some  $x, y \in \mathbb{Z}$ . Since  $d \mid a$  and  $d \mid b$ , the fundamental rules of divisibility imply that  $d \mid d'$ , and in particular  $d \leq d'$ .

We show that we also have  $d' \leq d$ . To this end, we divide  $a$  and  $b$  by  $d'$  with remainder, that is, we write

$$a = q_1d' + r_1, \quad b = q_2d' + r_2,$$

with  $0 \leq r_1, r_2 < d'$ . Using  $d' = ax + by$  we obtain

$$\begin{aligned} r_1 &= a - q_1d' = a - q_1(ax + by) = ax_1 + by_1 \\ r_2 &= b - q_2d' = b - q_2(ax + by) = ax_2 + by_2 \end{aligned}$$

with integers  $x_1, y_1, x_2, y_2 \in \mathbb{Z}$ . If we had  $r_1 > 0$ , then  $r_1 \in M$  would be smaller than  $d'$ , which is a contradiction to the minimality of  $d'$ . We find  $r_1 = 0$  and, analogously,  $r_2 = 0$ . Hence  $d'$  divides both  $a$  and  $b$ , which implies  $d' \leq d$ .  $\square$

From Bézout's Lemma we obtain a handy criterion to check whether two numbers are coprime.

**Corollary 1.2.5.** *We have  $\gcd(a, b) = 1$  if and only if there exist  $x, y \in \mathbb{Z}$  with  $ax + by = 1$ .*

We can now prove Euclid's Lemma (in a slightly stronger version), and thereby finish the proof of the Fundamental Theorem of Arithmetic.



## 1.2 The greatest common divisor and the Euclidean algorithm

**Lemma 1.2.6** (Euclid). *Let  $a, b, n \in \mathbb{N}$ . If  $n \mid ab$  and  $\gcd(a, n) = 1$ , then  $n \mid b$ .*

*Proof.* If  $\gcd(a, n) = 1$ , then by Bézout's Lemma there exist  $x, y \in \mathbb{Z}$  such that  $nx + ay = 1$ . Since  $n \mid ab$  there is  $k \in \mathbb{Z}$  such that  $nk = ab$ , so we obtain

$$b = b \cdot 1 = b(nx + ay) = bnx + aby = bnx + kny = n(bx + ky),$$

which means that  $n \mid b$ . □

In order to compute the gcd, we may use the *Euclidean algorithm*.

**Theorem 1.2.7** (Euclidean Algorithm). *Let  $a, b \in \mathbb{Z}$  not both 0 and assume that  $|a| \geq |b|$  (otherwise interchange  $a$  and  $b$ ). Define non-negative integers  $r_k, q_k$  by repeatedly performing division with remainder as follows:*

$$\begin{aligned} a &= q_1 b + r_1 & (0 \leq r_1 < |b|), \\ b &= q_2 r_1 + r_2 & (0 \leq r_2 < r_1), \\ r_1 &= q_3 r_2 + r_3 & (0 \leq r_3 < r_2), \\ &\vdots \\ r_{n-2} &= q_n r_{n-1} + r_n & (0 \leq r_n < r_{n-1}) \\ r_{n-1} &= q_{n+1} r_n + \underbrace{r_{n+1}}_{=0}, \end{aligned}$$

until the remainder  $r_{n+1}$  becomes 0. Then  $r_n = \gcd(a, b)$ .

Moreover, define integers  $x_k, y_k$  by

$$\begin{aligned} x_0 &= 0, & y_0 &= 1, \\ x_1 &= 1, & y_1 &= -q_1, \\ x_k &= x_{k-2} - q_k x_{k-1}, & y_k &= y_{k-2} - q_k y_{k-1}. \end{aligned}$$

Put  $x = x_n$  and  $y = y_n$ . Then we have

$$\gcd(a, b) = ax + by.$$

*Proof.* Since the remainders  $r_i$  are non-negative integers which get smaller in every step, the algorithm terminates after finitely many steps. The fact that we indeed have  $r_n = \gcd(a, b)$  follows by induction from the simple fact that  $a = bq + r$  implies  $\gcd(a, b) = \gcd(b, r)$ . We leave the details as an exercise.

We prove by induction that

$$ax_k + by_k = r_k$$

for  $k \geq 1$ . For  $k = 1$  we have  $x_1 = 1$  and  $y_1 = -q_1$ , and we indeed have

$$ax_1 + by_1 = a - q_1 b = r_1.$$

Let us assume that  $ax_j + by_j = r_j$  holds for all  $j \leq k$ , for some fixed  $k$ . We compute

$$\begin{aligned} ax_{k+1} + by_{k+1} &= a(x_{k-1} - q_{k+1}x_k) + b(y_{k-1} - q_{k+1}y_k) \\ &= (ax_{k-1} + by_{k-1}) - q_{k+1}(ax_k + by_k) \\ &= r_{k-1} - q_{k+1}r_k \\ &= r_{k+1}, \end{aligned}$$

by induction hypothesis and the definition  $r_{k+1} = q_{k+1}r_k + r_{k-1}$ . This finishes the proof. □

## 1 Primes and Divisibility

**Example 1.2.8.** Let us compute  $\gcd(43, 12)$  with the Euclidean Algorithm. We perform division with remainder:

$$\begin{aligned}43 &= 3 \cdot 12 + 7 & (a &= q_1b + r_1) \\12 &= 1 \cdot 7 + 5 & (b &= q_2r_1 + r_2) \\7 &= 1 \cdot 5 + 2 & (r_1 &= q_3r_2 + r_3) \\5 &= 2 \cdot 2 + 1 & (r_2 &= q_4r_3 + r_4) \\2 &= 2 \cdot 1 + 0 & (r_3 &= q_5r_4 + r_5)\end{aligned}$$

The algorithm finishes since the last remainder  $r_5$  is 0. The second-to-last remainder  $r_4$  is 1, so we obtain  $\gcd(43, 12) = r_4 = 1$ .

Let us compute  $x, y$  with  $43x + 12y = 1$ . We need to compute  $x_4, y_4$ . We have

$$\begin{aligned}x_0 &= 0, & y_0 &= 1, \\x_1 &= 1, & y_1 &= -3, \\x_2 &= 0 - 1 \cdot 1 = -1, & y_2 &= 1 - 1 \cdot (-3) = 4, \\x_3 &= 1 - 1 \cdot (-1) = 2, & y_3 &= -3 - 1 \cdot 4 = -7, \\x_4 &= -1 - 2 \cdot 2 = -5, & y_4 &= 4 - 2 \cdot (-7) = 18.\end{aligned}$$

Hence we take  $x = x_4 = -5$  and  $y = y_4 = 18$ , and indeed find that  $43 \cdot (-5) + 12 \cdot 18 = 1$ .

## 1.3 Distribution of primes

### 1.3.1 Bertrand's Postulate

It is easy to see that the gaps between primes can be arbitrary large.

**Proposition 1.3.1.** *For every  $n \in \mathbb{N}$  there exist  $n$  consecutive natural numbers which are all not prime.*

*Proof.* We consider the  $n$  consecutive numbers

$$(n+1)! + 2, (n+1)! + 3, \dots, (n+1)! + n + 1,$$

where  $n! = 1 \cdot 2 \cdot 3 \cdots (n-1) \cdot n$  denotes the factorial. For every  $2 \leq k \leq n+1$  we have  $k \mid (n+1)! + k$ , so none of the above  $n$  natural numbers is prime.  $\square$

However, such large prime gaps can only occur between very large numbers, as the following important result shows.

**Theorem 1.3.2** (Bertrand's Postulate). *For every  $n \in \mathbb{N}$  there is a prime  $p$  with  $n < p \leq 2n$ .*

This theorem was conjectured by Bertrand in 1845 and first proved by Chebychev in 1850. We will discuss an elementary proof by Erdős from 1932. The proof closely follows the exposition in "Proofs from THE BOOK" by Martin Aigner and Günter M. Ziegler, and takes up the rest of this section.

The key is to analyze the possible prime factors and the growth of the binomial coefficient

$$\binom{2n}{n} = \frac{(2n)!}{(n!)^2}.$$

Here  $\binom{n}{k} = \frac{n!}{k!(n-k)!}$ . Notice that each prime factor  $p$  of  $\binom{2n}{n}$  satisfies  $p \leq 2n$ . In particular, the primes  $p$  with  $n < p \leq 2n$  are the “largest” possible prime factors of  $\binom{2n}{n}$ . Hence, if we assume that no such prime exists (contradicting Bertrand’s postulate), we might expect that  $\binom{2n}{n}$  is rather “small”. On the other hand, there is an elementary lower bound

$$\binom{2n}{n} \geq \frac{4^n}{2n},$$

which tells us that  $\binom{2n}{n}$  is rather “large”, yielding a contradiction. This lower bound will be discussed in the exercises.

Hence, the main point in the proof of Bertrand’s Postulate is to find a good *upper* bound for  $\binom{2n}{n}$ . To derive this upper bound, we study the prime factors of  $\binom{2n}{n}$ . We divide these primes into the “small” prime factors  $p \leq \sqrt{2n}$  and the “large” prime factors  $p > \sqrt{2n}$ . The small prime factors will not bother us too much since they do not contribute much to the growth of  $\binom{2n}{n}$ , but we need to investigate the large prime factors in some detail. We have the following result about the large prime factors:

**Lemma 1.3.3.**

1. Let  $p > \sqrt{2n}$ . Then  $p$  occurs at most once in the prime factorization of  $\binom{2n}{n}$ .
2. Let  $\frac{2}{3}n < p \leq n$ . Then  $p$  does not divide  $\binom{2n}{n}$ .

For the proof we need the following result, which we leave as an exercise.

**Theorem 1.3.4** (Legendre). *The multiplicity of  $p$  in  $n!$  is*

$$\sum_{k=1}^{\infty} \left\lfloor \frac{n}{p^k} \right\rfloor.$$

Here  $\lfloor x \rfloor$  denotes the largest integer  $\leq x$ , and the sum is actually finite.

*Proof of Lemma 1.3.3.* We only prove the first part, and leave the second item as an exercise.

By Legendre’s Theorem, the multiplicity of  $p$  in  $\binom{2n}{n} = \frac{(2n)!}{(n!)^2}$  is given by

$$\sum_{k=1}^{\infty} \left( \left\lfloor \frac{2n}{p^k} \right\rfloor - 2 \left\lfloor \frac{n}{p^k} \right\rfloor \right).$$

Here, each summand is at most 1, since it satisfies

$$\underbrace{\left\lfloor \frac{2n}{p^k} \right\rfloor}_{\leq \frac{2n}{p^k}} - 2 \underbrace{\left\lfloor \frac{n}{p^k} \right\rfloor}_{> \frac{n}{p^k} - 1} < \frac{2n}{p^k} - 2 \left( \frac{n}{p^k} - 1 \right) = 2,$$

and it is an integer. Furthermore, the summands vanish whenever  $p^k > 2n$ . Thus, the multiplicity of  $p$  in  $\binom{2n}{n}$  can be estimated by

$$\sum_{k=1}^{\infty} \left( \left\lfloor \frac{2n}{p^k} \right\rfloor - 2 \left\lfloor \frac{n}{p^k} \right\rfloor \right) \leq \max\{r : p^r \leq 2n\}.$$

This means that the largest power of  $p$  that divides  $\binom{2n}{n}$  is not larger than  $2n$ . In particular, primes  $p > \sqrt{2n}$  appear at most once in  $\binom{2n}{n}$ . □

## 1 Primes and Divisibility

The above lemma implies that the possible primes occurring in  $\binom{2n}{n}$  are

1. some “small” primes  $p \leq \sqrt{2n}$ , possibly with multiplicity  $> 1$ .
2. some “large” primes  $\sqrt{2n} < p \leq \frac{2}{3}n$ , but each at most once,
3. no primes  $\frac{2}{3}n < p \leq n$ ,
4. some “large” primes  $n < p \leq 2n$ , but each at most once.

**Example 1.3.5.** For  $n = 25$  we have

$$\binom{50}{25} = 126410606437752 = 2^3 \cdot 3^2 \cdot 7^2 \cdot 13 \cdot 29 \cdot 31 \cdot 37 \cdot 41 \cdot 43 \cdot 47.$$

Note that  $\sqrt{2 \cdot 25} \sim 7.07$  and  $\frac{2}{3} \cdot 25 \sim 16.67$ . We see that the small primes  $p \leq 7$  indeed appear to higher powers, but the primes  $p > 7$  appear at most once. As Lemma 1.3.3 predicts, there are no primes between 16.67 and 25.

In the proof of the last lemma we have also seen that the largest power of a prime dividing  $\binom{2n}{n}$  is bounded by  $2n$ . Hence, as a preliminary upper bound we get

$$\binom{2n}{n} \leq \prod_{p \leq \sqrt{2n}} 2n \cdot \prod_{\sqrt{2n} < p \leq \frac{2}{3}n} p \cdot \prod_{n < p \leq 2n} p, \quad (1.3.1)$$

where we estimated the small prime powers appearing in  $\binom{2n}{n}$  by  $2n$ . We will further estimate the product over the small primes trivially by  $(2n)^{\sqrt{2n}}$ , but for the second product, we need to be more careful. Here we will use the following estimate.

**Lemma 1.3.6.** For all  $x \in \mathbb{R}$  with  $x \geq 2$  we have

$$\prod_{p \leq x} p \leq 4^{x-1}.$$

*Proof.* Note that, if  $q = \lfloor x \rfloor$  is the largest integer  $\leq x$ , then

$$\prod_{p \leq x} p = \prod_{p \leq q} p \quad \text{and} \quad 4^{q-1} \leq 4^{x-1},$$

so it suffices to prove the lemma in the case that  $x = q$  is a positive integer. In this case, we can prove the claim by induction.

For  $q = 2$  we get  $2 \leq 4$ , which is true. Let us now suppose that the lemma is true for all natural numbers  $x$  less than some fixed *odd* integer  $q = 2m + 1$ . Then we split the product as

$$\prod_{p \leq 2m+1} p = \prod_{p \leq m+1} p \cdot \prod_{m+1 < p \leq 2m+1} p \leq 4^m \binom{2m+1}{m} \leq 4^m 2^{2m} = 4^{2m}.$$

Here we used the induction hypothesis

$$\prod_{p \leq m+1} p \leq 4^m,$$

the inequality

$$\prod_{m+1 < p \leq 2m+1} p \leq \binom{2m+1}{m},$$

(exercise!) and, finally, the fact that

$$\binom{2m+1}{m} \leq 2^{2m}$$

(exercise!). This finishes the induction step for *odd* integers  $q$ . But if  $q = 2m > 2$  is *even*, then  $q$  is not a prime, so we have

$$\prod_{p \leq q} p = \prod_{p \leq q-1} p \leq 4^{(q-1)-1} < 4^{q-1},$$

where we used the induction step for odd  $q$  proved above. This finishes the proof.  $\square$

Using Lemma 1.3.6 we can further estimate (1.3.2) to obtain our final upper bound:

**Corollary 1.3.7.** *We have*

$$\binom{2n}{n} \leq (2n)^{\sqrt{2n}} \cdot 4^{\frac{2}{3}n} \cdot \prod_{n < p \leq 2n} p. \quad (1.3.2)$$

With this upper bound, we can now prove Bertrand's Postulate.

*Proof of Theorem 1.3.2.* Let us assume that there exists some  $n$  such that Bertrand's Postulate is violated, that is, there is no prime  $p$  with  $n < p \leq 2n$ . Then the last product in the upper bound (1.3.2) is empty, that is, it equals 1. Thus, we obtain the upper bound

$$\binom{2n}{n} \leq (2n)^{\sqrt{2n}} \cdot 4^{\frac{2}{3}n}.$$

Combining this with the elementary lower bound

$$\binom{2n}{n} \geq \frac{4^n}{2n},$$

we obtain the estimate

$$4^{\frac{1}{3}n} \leq (2n)^{1+\sqrt{2n}}.$$

However, it is easy to check that this is false for large  $n$ , e.g.  $n > 4000$ . Hence, for  $n > 4000$  Bertrand's Postulate must be true. So it remains to check it for  $n \leq 4000$ . To this end, notice that each of the primes in the sequence

$$2, 3, 5, 7, 13, 23, 43, 83, 163, 317, 631, 1259, 2503, 4001$$

is less than twice as big as its predecessor, so every interval of the form  $(n, 2n]$  with  $n \leq 4000$  contains at least one of the above primes. This finishes the proof of Bertrand's Postulate.  $\square$

## 1 Primes and Divisibility

We remark that Bertrand's Postulate is not optimal if  $n$  is large. For example, in 1952 Nagura proved that for  $n \geq 25$  there is always a prime between  $n$  and  $(1 + \frac{1}{5})n$ , and there are similar (stronger) bounds for even larger  $n$ . For comparison, Bertrand's Postulate predicts that the next prime after 1000 has to appear somewhere before 2000, whereas Nagura's result says that the next prime appears before 1200. In fact, it follows from the Prime Number Theorem (that we discuss in the next section) that for any  $\varepsilon > 0$  there is a  $n_0 > 0$  such that for all  $n > n_0$  there is a prime  $p$  such that  $n < p < (1 + \varepsilon)n$ .

Moreover, Bertrand's Postulate states that there is always *at least one* prime between  $n$  and  $2n$ , but there are usually many primes in such an interval. More precisely, Ramanujan in 1919 showed that the number of primes between  $n$  and  $2n$  goes to infinity as  $n \rightarrow \infty$ .

There are several similar results and conjectures about primes in certain intervals. For example, Legendre's conjecture states that there should always be a prime between  $n^2$  and  $(n + 1)^2$ , which is an open problem.

### 1.3.2 The Prime Number Theorem

For real  $x > 0$  we consider the *prime counting function*

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

Using tools from complex analysis, one can show the *Prime Number Theorem*.

**Theorem 1.3.8** (Hadamard / de la Vallée Poussin). *We have*

$$\lim_{x \rightarrow \infty} \frac{\pi(x)}{\frac{x}{\log(x)}} = 1,$$

where  $\log(x) = \ln(x)$  denotes the natural logarithm (to base  $e$ ).

This result was conjectured independently by Gauss and Legendre in 1792 and 1798, respectively, and, again independently, proved by de la Vallée Poussin and Hadamard in 1896.

The goal of this section is to derive the weaker estimates

$$\frac{2}{3} \frac{x}{\log(x)} \leq \pi(x) \leq 2 \frac{x}{\log(x)},$$

using elementary methods. For the lower bound, we use the following lemma.

**Lemma 1.3.9.** *Let  $n, k \in \mathbb{N}$  with  $k \leq n$ . Then we have the estimate*

$$\binom{n}{k} \leq n^{\pi(n)}.$$

*Proof.* Recall that  $\nu_p(n)$  denotes the multiplicity of the prime  $p$  in the prime factorization of  $n$ . We first show that

$$\nu_p \left( \binom{n}{k} \right) \leq \frac{\log(n)}{\log(p)}. \tag{1.3.3}$$

Indeed, using Legendre's Theorem 1.3.4, we compute

$$\begin{aligned} \nu_p \left( \binom{n}{k} \right) &= \nu_p \left( \frac{n!}{k!(n-k)!} \right) \\ &= \nu_p(n!) - \nu_p(k!) - \nu_p((n-k)!) \\ &= \sum_{j=1}^{\infty} \left( \left\lfloor \frac{n}{p^j} \right\rfloor - \left\lfloor \frac{k}{p^j} \right\rfloor - \left\lfloor \frac{n-k}{p^j} \right\rfloor \right) \\ &\leq \frac{\log(n)}{\log(p)}, \end{aligned}$$

since for  $j > \frac{\log(n)}{\log(p)}$  every term in the sum is 0, and for  $j \leq \frac{\log(n)}{\log(p)}$  every summand is 0 or 1. Multiplying by  $\log(p)$  on both sides of (1.3.3) and applying the exponential  $e^x$ , we obtain

$$p^{\nu_p \left( \binom{n}{k} \right)} \leq n.$$

Now for each prime  $p \mid \binom{n}{k}$  we have  $p \leq n$ , hence

$$\binom{n}{k} = \prod_{p \leq n} p^{\nu_p \left( \binom{n}{k} \right)} \leq \prod_{p \leq n} n = n^{\pi(n)}.$$

This finishes the proof. □

**Proposition 1.3.10.** *For  $n \in \mathbb{N}$  with  $n \geq 2$  we have*

$$\pi(n) > \frac{2}{3} \frac{n}{\log(n)}.$$

*Proof.* Using  $\binom{n}{0} = \binom{n}{n} = 1$  and the last lemma, we have

$$2^n = (1+1)^n = \sum_{k=0}^n \binom{n}{k} \leq 1 + (n-1)n^{\pi(n)} + 1 \leq n^{\pi(n)+1},$$

where we used  $2 \leq n^{\pi(n)}$  for  $n \geq 2$  in the last step. Taking logarithms we get

$$\pi(n) \geq \log(2) \frac{n}{\log(n)} - 1.$$

It is easy to check that for  $n > 200$  we have

$$\log(2) \frac{n}{\log(n)} - 1 > \frac{2}{3} \frac{n}{\log(n)},$$

and the remaining cases can be checked by hand. □

For the upper bound for  $\pi(n)$ , we need the following lemma.

**Lemma 1.3.11.** *For every  $n \in \mathbb{N}$  with  $n \geq 2$  we have*

$$\pi(2n) - \pi(n) < \log(4) \frac{n}{\log(n)}.$$

## 1 Primes and Divisibility

*Proof.* Note that  $\pi(2n) - \pi(n)$  is precisely the number of primes in the interval  $(n, 2n]$ . Each of these primes appears in the numerator of  $\binom{2n}{n} = \frac{(2n)!}{(n!)^2}$ , but not in the denominator. Hence we can estimate

$$n^{\pi(2n) - \pi(n)} \leq \prod_{n < p \leq 2n} p \leq \binom{2n}{n} < 4^n,$$

where we used in the last step that  $\binom{2n}{n}$  appears in the sum  $\sum_{k=0}^{2n} \binom{2n}{k} = 2^{2n}$  and hence must be less than  $4^n$ . Taking logarithms gives the result.  $\square$

Now we obtain an upper bound for  $\pi(n)$ .

**Proposition 1.3.12.** *For every  $n \in \mathbb{N}$  with  $n \geq 2$  we have*

$$\pi(n) < 2 \frac{n}{\log(n)}.$$

*Proof.* We prove this by induction on  $n$ . For  $n \leq 256 = 2^8$  the claim can be proved directly. Suppose the claim is true for all  $k < n$  for some  $n$ .

First suppose that  $n = 2m > 2^8$  is even. Using the last lemma and the induction hypothesis we find

$$\pi(n) = \pi(2m) < \pi(m) + \log(4) \frac{m}{\log(m)} < 2 \frac{m}{\log(m)} + \log(4) \frac{m}{\log(m)} = (1 + \log(2)) \frac{n}{\log(n) - \log(2)}.$$

But for  $n \geq 2^8$  we have

$$\frac{1}{\log(n) - \log(2)} \leq \frac{8}{7} \frac{1}{\log(n)},$$

and using  $(1 + \log(2)) \frac{8}{7} < 2$  we obtain  $\pi(n) < 2 \frac{n}{\log(n)}$ .

If  $n = 2m + 1 > 2^8$  is odd, then we obtain as above

$$\begin{aligned} \pi(n) &= \pi(2m + 1) \leq \pi(2m) + 1 < (1 + \log(2)) \frac{8}{7} \frac{2m}{\log(2m)} + 1 \\ &< (1 + \log(2)) \frac{8}{7} \frac{n}{\log(n)} + 1 = \left( (1 + \log(2)) \frac{8}{7} + \frac{\log(n)}{n} \right) \frac{n}{\log(n)}. \end{aligned}$$

Using that  $x \mapsto \frac{\log(x)}{x}$  is monotonically decreasing for  $x \geq e$ , we obtain for  $n > 2^8$  that

$$(1 + \log(2)) \frac{8}{7} + \frac{\log(n)}{n} \leq (1 + \log(2)) \frac{8}{7} + \frac{8 \log(2)}{256} < 2.$$

This finishes the proof.  $\square$

We finish this section with some famous conjectures about primes numbers.

### The twin prime conjecture

Two primes which differ by 2 are called *twin primes*, for example (3, 5), (5, 7), (11, 13), (101, 103).

**Conjecture 1.3.13.** *There are infinitely many twin primes.*

A major breakthrough was obtain by Yitang Zhang in 2013 when he proved that there are infinitely many pairs of primes which differ by at most 70 million. Shortly after, Terence Tao, James Maynard, and others improved the gap to 246. Assuming some other strong conjectures the gap can be improved to 6. However, it seems that a proof of the twin prime conjecture will require new methods.



**Goldbach's conjecture**

The following was conjectured by Christian Goldbach in 1642.

**Conjecture 1.3.14** (Strong form of Goldbach's conjecture). *Every even number  $> 2$  can be written as a sum of two primes.*

The best proved result in this direction is due to Olivier Ramaré, who showed in 1995 that every even number  $> 2$  can be written as a sum of at most six primes.

**Conjecture 1.3.15** (Weak form of Goldbach's conjecture). *Every odd number  $> 5$  can be written as a sum of three primes.*

The strong form would imply the weak form: if  $n > 5$  is odd, then  $n - 3 > 2$  is even, so  $n - 3 = p + q$  gives  $n = p + q + 3$ .

In 2012, Terence Tao proved that each odd number  $> 1$  can be written as a sum of five primes. In 2013, Harald Helfgott announced a proof of the weak Goldbach conjecture, which is now widely accepted to be correct.

**Primes of special forms**

By Dirichlet's Theorem on primes in arithmetic progressions, we know that there are infinitely many primes of the form  $an + b$ , if  $\gcd(a, b) = 1$ . One can ask if there are other natural families of numbers which contain infinitely many primes. For example:

**Conjecture 1.3.16** (Landau 1912). *There are infinitely many primes of the form  $n^2 + 1$ .*

Two important families are given by the *Mersenne numbers*  $M_n = 2^n - 1$  and the *Fermat numbers*  $F_n = 2^{2^n} + 1$ .

There are currently 51 known Mersenne primes<sup>1</sup>, and it is conjectured that there are infinitely many of them. However, there are only five known Fermat primes,  $F_0 = 3, F_1 = 5, F_2 = 17, F_3 = 257, F_4 = 65537$ , and there are heuristics that say that there should be no others.

The largest known primes are all Mersenne primes, the current record being  $2^{82,589,933} - 1$ . The reason why Mersenne primes are used to construct large primes is that the *Lucas-Lehmer test* provides an efficient algorithm to check whether a given Mersenne number is prime.

---

<sup>1</sup>See <https://www.mersenne.org> for the current status on the search for Mersenne primes.



## 2 Number-theoretic functions

### 2.1 Multiplicative number-theoretic functions

**Definition 2.1.1.** A *number-theoretic function*<sup>1</sup> is a function

$$f : \mathbb{N} \rightarrow \mathbb{C}.$$

We call  $f$  *multiplicative* if

$$f(mn) = f(m)f(n) \quad \text{for all coprime } m, n \in \mathbb{N},$$

and *completely multiplicative* if this identity holds for *all*  $m, n \in \mathbb{N}$ .

We collect some important examples of number-theoretic functions.

**Example 2.1.2.** 1. The identity function  $\text{id}(n) = n$  and

$$e(n) = \begin{cases} 1, & \text{if } n = 1, \\ 0, & \text{otherwise,} \end{cases}$$

are completely multiplicative number-theoretic functions.

2. The only constant multiplicative number-theoretic functions are the constant 0-function and the constant 1-function  $\mathbf{1}(n) = 1$ .

3. Let  $k \in \mathbb{Z}$ . The *divisor sum*

$$\sigma_k(n) = \sum_{d|n} d^k$$

is multiplicative, but not completely multiplicative. Here the sum  $\sum_{d|n}$  runs over all positive divisors of  $n \in \mathbb{N}$ . We will show in the exercises that  $\sigma_k(n)$  is multiplicative.

We also write

$$\tau(n) = \sigma_0(n) = \#\{d \in \mathbb{N} : d | n\}$$

for the *number-of-divisors function*, and

$$\sigma(n) = \sigma_1(n) = \sum_{d|n} d$$

for the *sum-of-divisors function*.

---

<sup>1</sup>Number-theoretic functions are also called *arithmetic functions*. Moreover, note that a number-theoretic function can also be viewed as a sequence  $(f(n))_{n \in \mathbb{N}}$  of complex numbers.

## 2 Number-theoretic functions

4. A very important number-theoretic function is *Euler's totient function*<sup>2</sup>

$$\varphi(n) = \#\{k \in \mathbb{N} : 1 \leq k \leq n \text{ and } \gcd(k, n) = 1\}.$$

We will show below that  $\varphi(n)$  is multiplicative. However, it is not completely multiplicative. Moreover, in the exercises we will derive the explicit formula

$$\varphi(n) = n \prod_{p|n} \left(1 - \frac{1}{p}\right),$$

where the product  $\prod_{p|n}$  runs over all prime divisors of  $n$ .

5. The *Möbius function*

$$\mu(n) = \begin{cases} (-1)^r, & \text{if } n \text{ is square-free, } n = p_1 \cdots p_r, \\ 0, & \text{if } n \text{ is not square-free,} \end{cases}$$

is multiplicative, but not completely multiplicative.

The following basic facts are easy to prove:

**Lemma 2.1.3.** *Let  $f, g : \mathbb{N} \rightarrow \mathbb{C}$  be multiplicative number-theoretic functions.*

1. *If  $f$  is not the constant 0-function, then  $f(1) = 1$ .*
2. *The function  $f \cdot g$  is multiplicative.*

A useful property of multiplicative function is that it suffices to study their values at prime powers.

**Proposition 2.1.4.** *Let  $f : \mathbb{N} \rightarrow \mathbb{C}$  be a number-theoretic function. Then the following are equivalent.*

1.  *$f$  is multiplicative.*
2. *For all  $n \in \mathbb{N}$  with prime factorization  $n = p_1^{\nu_1} \cdots p_r^{\nu_r}$  we have*

$$f(n) = f(p_1^{\nu_1}) \cdots f(p_r^{\nu_r}).$$

*In particular, a multiplicative number-theoretic function is uniquely determined by its values  $f(p^\nu)$  on prime powers  $p^\nu$  for  $\nu \in \mathbb{N}$ .*

*Proof.* The implication 1.  $\Rightarrow$  2. is clear by induction. For the converse implication, let  $m, n \in \mathbb{Z}$  with  $\gcd(m, n) = 1$ . Let

$$m = p_1^{\nu_1} \cdots p_r^{\nu_r}, \quad n = q_1^{\mu_1} \cdots q_s^{\mu_s},$$

be the prime factorizations of  $m, n$ . Since  $m, n$  are coprime, the primes  $p_1, \dots, p_r, q_1, \dots, q_s$  are pairwise different, so using 2. we obtain

$$f(mn) = f(p_1^{\nu_1} \cdots p_r^{\nu_r} \cdot q_1^{\mu_1} \cdots q_s^{\mu_s}) = f(p_1^{\nu_1}) \cdots f(p_r^{\nu_r}) \cdot f(q_1^{\mu_1}) \cdots f(q_s^{\mu_s}) = f(m) \cdot f(n),$$

hence  $f$  is multiplicative. □

---

<sup>2</sup>Some authors write  $\phi$  instead of  $\varphi$  for Euler's totient function.

**Definition 2.1.5.** For a number-theoretic function  $f$  we call

$$F(n) = \sum_{d|n} f(d)$$

the *summatory function* of  $f$ .

**Example 2.1.6.** 1. The number-of-divisors function  $\tau$  and the sum-of-divisors function  $\sigma$  are the summatory functions of  $\mathbf{1}$  and  $\text{id}$ , respectively. More generally,  $\sigma_k(n)$  is the summatory function of  $f(n) = n^k$ .

2. The summatory function of Euler's totient function  $\varphi$  is  $\text{id}$ . Explicitly, this means that

$$\sum_{d|n} \varphi(d) = n$$

for all  $n \in \mathbb{N}$ . This will be proved below.

3. The summatory function of the Möbius function is  $e$ . In other words,  $\mu = \widehat{\mathbf{1}}$  is the inverse of the constant 1-function with respect to convolution. In formulas, this means that

$$\sum_{d|n} \mu(d) = \begin{cases} 1, & \text{if } n = 1, \\ 0, & \text{otherwise.} \end{cases}$$

This will be proved in the exercises.

**Proposition 2.1.7.** *Let  $f$  be a number-theoretic function. Then  $f$  is multiplicative if and only if its summatory function is multiplicative.*

For the proof of the proposition, the following simple lemma is useful.

**Lemma 2.1.8.** *If  $m, n \in \mathbb{N}$  are coprime, then every divisor  $d \mid mn$  can be written uniquely as  $d = d_1 d_2$  with  $d_1 \mid m$  and  $d_2 \mid n$ .*

We will leave the proof of the proposition and the lemma as an exercise. The proposition will also follow from the results about Dirichlet convolution and Möbius inversion below.

## 2.2 Dirichlet convolution and Möbius inversion

**Definition 2.2.1.** For two number-theoretic functions  $f, g : \mathbb{N} \rightarrow \mathbb{C}$  we define their *Dirichlet convolution*  $f * g$  as the number-theoretic function

$$(f * g)(n) = \sum_{d|n} f(d)g\left(\frac{n}{d}\right).$$

**Example 2.2.2.** The Dirichlet convolution of a number-theoretic function  $f$  with the constant 1-function  $\mathbf{1}(n) = 1$  is just the summatory function of  $f$ ,

$$(f * \mathbf{1})(n) = \sum_{d|n} f(d).$$

## 2 Number-theoretic functions

The following lemma shows that the set of number-theoretic functions becomes a ring under pointwise addition and Dirichlet convolution, with neutral element

$$e(n) = \begin{cases} 1, & \text{if } n = 1, \\ 0, & \text{otherwise.} \end{cases}$$

**Lemma 2.2.3.** *Dirichlet convolution has the following properties.*

1. *Commutative:*  $f * g = g * f$
2. *Associative:*  $f * (g * h) = (f * g) * h$
3. *Neutral element:*  $f * e = e * f = f$
4. *Distributive:*  $(f + g) * h = f * h + g * h$ .

*Proof.* We prove 1. and 2. and leave the rest as an exercise. For 1., we have

$$(f * g)(n) = \sum_{d|n} f(d)g\left(\frac{n}{d}\right) = \sum_{d|n} f\left(\frac{n}{d}\right)g(d) = (g * f)(n),$$

where we just interchanged the order of summation. For 2., first note that we can write  $(f * g)(n) = \sum_{ab=n} f(a)g(b)$ . Hence, we have

$$(f * (g * h))(n) = \sum_{ab=n} f(a)(g * h)(b) = \sum_{ab=n} \sum_{cd=b} f(a)g(c)h(d) = \sum_{acd=n} f(a)g(c)h(d),$$

and

$$((f * g) * h)(n) = \sum_{ab=n} (f * g)(a)h(b) = \sum_{ab=n} \sum_{cd=a} f(c)g(d)h(b) = \sum_{bcd=n} f(c)g(d)h(b),$$

which is the same after renaming the variables. □

**Proposition 2.2.4.** *If  $f, g$  are multiplicative, then  $(f * g)$  is multiplicative, too.*

*Proof.* Using Lemma 2.1.8 we obtain for coprime  $m, n \in \mathbb{N}$ :

$$\begin{aligned} (f * g)(mn) &= \sum_{d|mn} f(d)g\left(\frac{mn}{d}\right) \\ &= \sum_{d_1|m, d_2|n} f(d_1d_2)g\left(\frac{m}{d_1} \cdot \frac{n}{d_2}\right) \\ &= \sum_{d_1|m, d_2|n} f(d_1)f(d_2)g\left(\frac{m}{d_1}\right)g\left(\frac{n}{d_2}\right) \\ &= \left(\sum_{d_1|m} f(d_1)g\left(\frac{m}{d_1}\right)\right) \left(\sum_{d_2|n} f(d_2)g\left(\frac{n}{d_2}\right)\right) \\ &= (f * g)(m) \cdot (f * g)(n). \end{aligned}$$

□

**Proposition 2.2.5.** *Let  $f$  be a number theoretic function with  $f(1) \neq 0$ .*

1. *There exists a unique number-theoretic function  $\widehat{f}$  with*

$$f * \widehat{f} = \widehat{f} * f = e,$$

*that is,  $\widehat{f}$  is the inverse of  $f$  with respect to convolution.*

2. *If  $f$  is multiplicative, then  $\widehat{f}$  is multiplicative.*

*Proof.* 1. We want to determine  $\widehat{f} : \mathbb{N} \rightarrow \mathbb{C}$  such that

$$(f * \widehat{f})(1) = 1 \quad \text{and} \quad (f * \widehat{f})(n) = 0 \quad \text{for all } n > 1.$$

Recall the definition

$$(f * \widehat{f})(n) = \sum_{d|n} f(d) \widehat{f}\left(\frac{n}{d}\right). \quad (2.2.1)$$

In particular, we must have  $\widehat{f}(1) = 1/f(1)$ , which explains the assumption  $f(1) \neq 0$ . Moreover, we can use (2.2.1) inductively to define  $\widehat{f}$ . Let us assume that  $n > 1$ , and that we already defined  $\widehat{f}(n')$  for all  $n' < n$ . Then, by writing

$$(f * \widehat{f})(n) = \sum_{d|n} f(d) \widehat{f}\left(\frac{n}{d}\right) = f(1) \widehat{f}(n) + \sum_{\substack{d|n \\ d \neq 1}} f(d) \widehat{f}\left(\frac{n}{d}\right)$$

we see that we can define  $\widehat{f}(n)$  by

$$\widehat{f}(n) = -\frac{1}{f(1)} \sum_{\substack{d|n \\ d \neq 1}} f(d) \widehat{f}\left(\frac{n}{d}\right).$$

Note that the sum on the right-hand side only involves  $\widehat{f}(n')$  for  $n' < n$ , which is already defined by induction hypothesis. This shows the existence of  $\widehat{f}$ . The uniqueness follows from the ring axioms, since inverse elements in a commutative ring are automatically unique if they exist.

2. The proof proceeds by induction on  $N = mn$ . The case  $N = 1$  follows from  $\widehat{f}(1) = \frac{1}{f(1)} = 1$  since  $f(1) = 1$  for any non-zero multiplicative function.

Let us now suppose that  $N > 1$  and  $\widehat{f}(ab) = \widehat{f}(a)\widehat{f}(b)$  for all  $a, b \in \mathbb{N}$  with  $\gcd(a, b) = 1$  and  $ab < N$ . Let  $N = mn$  be a factorization with coprime  $m, n \in \mathbb{N}$ . By Lemma 2.1.8, we have

$$e(mn) = \sum_{d|mn} f(d) \widehat{f}\left(\frac{mn}{d}\right) = \sum_{d_1|m, d_2|n} f(d_1) f(d_2) \widehat{f}\left(\frac{m}{d_1} \frac{n}{d_2}\right).$$

On the other hand, we also have

$$e(mn) = e(m)e(n) = \sum_{d_1|m, d_2|n} f(d_1) \widehat{f}\left(\frac{m}{d_1}\right) f(d_2) \widehat{f}\left(\frac{n}{d_2}\right),$$

## 2 Number-theoretic functions

so we get

$$\sum_{d_1|m, d_2|n} f(d_1)f(d_2)\widehat{f}\left(\frac{m}{d_1}\frac{n}{d_2}\right) = \sum_{d_1|m, d_2|n} f(d_1)\widehat{f}\left(\frac{m}{d_1}\right) f(d_2)\widehat{f}\left(\frac{n}{d_2}\right) \quad (2.2.2)$$

For  $d_1 \neq m$  and  $d_2 \neq n$  we have, by induction hypothesis, that  $\widehat{f}(d_1d_2) = \widehat{f}(d_1)\widehat{f}(d_2)$  since  $\gcd(d_1, d_2) = 1$  and  $d_1d_2 < mn = N$ . Hence, the corresponding summands in (2.2.2) cancel out, leaving only  $\widehat{f}(mn) = \widehat{f}(m)\widehat{f}(n)$ .  $\square$

**Remark 2.2.6.** In more algebraic terms, Proposition 2.2.5 says that the set of number theoretic functions  $f$  with  $f(1) \neq 0$  is a group with respect to convolution, and the non-zero multiplicative number-theoretic functions form a subgroup.

The main result of this section is the *Möbius inversion formula*, which shows that a number-theoretic function can be recovered from its summatory function.

**Theorem 2.2.7** (Möbius inversion). *Let  $f$  and  $F$  be number-theoretic functions. The following two properties are equivalent.*

1.  $F = f * \mathbf{1}$  is the summatory function of  $f$ , that is,  $F(n) = \sum_{d|n} f(d)$  for all  $n \in \mathbb{N}$ .
2.  $f = F * \mu$  is the convolution of  $F$  and  $\mu$ , that is,  $f(n) = \sum_{d|n} F(d)\mu(n/d)$  for all  $n \in \mathbb{N}$ .

*Proof.* 1.  $\Rightarrow$  2. Since  $e$  is the summatory function of  $\mu$ , that is,  $e = \mathbf{1} * \mu$ , we have

$$f = f * e = f * (\mathbf{1} * \mu) = (f * \mathbf{1}) * \mu = F * \mu.$$

2.  $\Rightarrow$  1. Similarly as above, we can compute

$$F = F * e = F * (\mu * \mathbf{1}) = (F * \mu) * \mathbf{1} = f * \mathbf{1}.$$

$\square$

Since Dirichlet convolution preserves multiplicative functions, we obtain the following useful result.

**Corollary 2.2.8.** *Let  $f$  be a number-theoretic function and  $F = f * \mathbf{1}$  its summatory function. Then  $f$  is multiplicative if and only if  $F$  is multiplicative.*

As an application, we show that Euler's totient function  $\varphi$  is multiplicative. To this end, we compute its summatory function.

**Proposition 2.2.9.** *The summatory function of  $\varphi$  is id, that is, for all  $n \in \mathbb{N}$  we have*

$$\sum_{d|n} \varphi(d) = n.$$

*Proof.* We write the set  $\underline{n} = \{1, \dots, n\}$  as a disjoint union

$$\underline{n} = \bigcup_{d|n} T_d(n)$$

with

$$T_d(n) = \{m \in \underline{n} : \gcd(m, n) = d\}.$$



Then we have

$$n = \sum_{d|n} \#T_d(n).$$

Moreover, since  $\gcd(m, n) = d$  is equivalent to  $\gcd(m/d, n/d) = 1$ , we have

$$\#T_d(n) = \varphi(n/d)$$

for every  $d \mid n$ . We find

$$n = \sum_{d|n} \#T_d(n) = \sum_{d|n} \varphi(n/d) = \sum_{d|n} \varphi(d),$$

where we just reordered the summation in the last step. This finishes the proof.  $\square$

Since the summatory function of  $\varphi$  is id, which is multiplicative, we obtain:

**Corollary 2.2.10.** *Euler's totient function  $\varphi(n)$  is multiplicative.*

## 2.3 Perfect and amicable numbers

As an application of multiplicative number-theoretic functions, we study perfect and amicable numbers.

**Definition 2.3.1.** A natural number  $n \in \mathbb{N}$  is called *perfect* if it is equal to the sum of its proper positive divisors (that is,  $\sigma(n) = 2n$ ).

For example, 6 is perfect since  $6 = 1 + 2 + 3$  (or  $\sigma(6) = 1 + 2 + 3 + 6 = 12 = 2 \cdot 6$ ). The following criterion gives a classification of the *even* perfect numbers.

**Theorem 2.3.2** (Euclid, Euler). *An even natural number  $n \in \mathbb{N}$  is perfect if and only if*

$$n = 2^{m-1}(2^m - 1) \quad \text{and} \quad 2^m - 1 \text{ is prime}$$

for some  $m \in \mathbb{N}$ .

*Proof.* The proof that each number of the form  $2^{m-1}(2^m - 1)$ , with  $2^m - 1$  prime, is perfect, is easy and will be discussed in the exercises.

Conversely, let  $n$  be perfect, that is,  $\sigma(n) = 2n$ . We write  $n = 2^k b$  with  $b \in \mathbb{N}$  odd and  $k \geq 1$ . Then

$$\sigma(n) = \sigma(2^k)\sigma(b) = (2^{k+1} - 1)\sigma(b).$$

Since  $n$  is perfect, we also have  $\sigma(n) = 2n = 2^{k+1}b$ , which implies

$$\sigma(b) = \frac{\sigma(n)}{2^{k+1} - 1} = \frac{2^{k+1}}{2^{k+1} - 1}b = b + c$$

with

$$c = \frac{b}{2^{k+1} - 1} > 0.$$

Since  $\sigma(b) \in \mathbb{N}$  we have  $c = \sigma(b) - b \in \mathbb{N}$ , hence  $c \mid b$ . Thus  $c$  appears in the sum  $\sigma(b)$ . Now  $\sigma(b) = b + c$  implies that  $c = 1$  and  $b = 2^{k+1} - 1$  must be prime.  $\square$

## 2 Number-theoretic functions

In order to find even perfect numbers, we need to check whether  $2^m - 1$  is prime. A necessary condition is that  $m$  itself must be prime.

**Lemma 2.3.3.** *Let  $m \in \mathbb{N}$ . If  $2^m - 1$  is prime, then  $m$  is also prime.*

This will also be proved in the exercises. The lemma leads to the following definition.

**Definition 2.3.4.** For a prime  $p$  we let  $M_p = 2^p - 1$  the  $p$ -th *Mersenne number*. If  $M_p$  is prime, we call it a *Mersenne prime*.

Hence, searching for even perfect numbers is the same as searching for Mersenne primes. It is not known if there are infinitely many Mersenne primes or (equivalently) infinitely many even perfect numbers, but it is conjectured.

Although we have a nice classification of the even perfect numbers, we do not know if any odd perfect numbers exist. However, odd perfect numbers, if they existed, would need to satisfy several strong conditions, which makes their existence rather unlikely. For example, an odd perfect number would have to be bigger than  $10^{1500}$ , and have at least 9 different prime divisors, the largest prime factor being bigger than  $10^8$ . We confine ourselves with two simpler results.

**Proposition 2.3.5.** *Let  $n$  be an odd perfect number. Then  $n$  has at least 3 different prime factors.*

*Proof.* Suppose that  $n = p^\nu q^\mu$  is a product of only two odd prime powers. We show that  $\sigma(n) < 2n$ . We can assume that  $p < q$  and hence  $p \geq 3$  and  $q \geq 5$ . We can now estimate

$$\begin{aligned} \sigma(n) &= \sigma(p^\nu)\sigma(q^\mu) \\ &= (1 + p + p^2 + \cdots + p^{\nu-1} + p^\nu) \cdot (1 + q + q^2 + \cdots + q^{\mu-1} + q^\mu) \\ &= \left(\frac{p^\nu - 1}{p - 1} + p^\nu\right) \cdot \left(\frac{q^\mu - 1}{q - 1} + q^\mu\right) \\ &\leq \left(\frac{p^\nu}{2} + p^\nu\right) \cdot \left(\frac{q^\mu}{4} + q^\mu\right) = \frac{15}{8}p^\nu q^\mu < 2n, \end{aligned}$$

which shows that  $n$  cannot be a perfect number. □

**Proposition 2.3.6** (Euler). *An odd perfect number (if it exists) must be of the form*

$$n = p^{2m+1}Q^2,$$

where  $p$  is an odd prime and  $Q$  is an odd natural number with  $\gcd(p, Q) = 1$ .

*Proof.* Suppose that  $n$  is an odd perfect number, and let  $n = p_1^{\nu_1} \cdots p_r^{\nu_r}$  be the prime factorization of  $n$ , with distinct odd primes  $p_j$ . Since  $n$  is perfect, we have  $\sigma(n) = 2n$ . On the other hand, since  $\sigma(n)$  is multiplicative, we have

$$2n = \sigma(n) = \sigma(p_1^{\nu_1}) \cdots \sigma(p_r^{\nu_r}) = (1 + p_1 + p_1^2 + \cdots + p_1^{\nu_1}) \cdots (1 + p_r + p_r^2 + \cdots + p_r^{\nu_r}).$$

Since  $n$  is odd, exactly one of the sums on the right is even, and the others are odd. This implies that precisely one of the exponents  $\nu_j$  is odd, and all others are even. Letting  $\nu_j = 2m + 1$  and  $p = p_j$ , and  $Q^2$  the product of the remaining prime powers of  $n$ , we obtain that  $n = p^{2m+1}Q^2$  with  $p, Q$  odd and  $\gcd(p, Q) = 1$ . □

**Definition 2.3.7.** Two natural numbers  $m, n \in \mathbb{N}$  are called *amicable* if  $m$  is the sum of the proper positive divisors of  $n$ , and  $n$  is the sum of the proper positive divisors of  $m$  (or, equivalently, if  $\sigma(m) = \sigma(n) = n + m$ ).

The first pair of amicable numbers is  $(220, 284)$  with  $\sigma(220) = \sigma(284) = 504$ . A possible method to find amicable numbers was given in the 9th century by Thabit.

**Theorem 2.3.8** (Thabit's rule). *Let  $k \in \mathbb{N}$  with  $k \geq 2$ . If the three numbers*

$$\begin{aligned} T_k &= 3 \cdot 2^k - 1 \\ T_{k-1} &= 3 \cdot 2^{k-1} - 1 \\ R_k &= 9 \cdot 2^{2k-1} - 1 \end{aligned}$$

*are all prime, then the numbers*

$$m = 2^k T_k T_{k-1} \quad \text{and} \quad n = 2^k R_k$$

*are amicable.*

The proof is a simple exercise, which we leave to the reader. Unfortunately, the theorem only seems to yield three pairs of amicable numbers, for  $k = 2, 4,$  and  $7$ . However, there exist many variants of Thabit's rule which yield many new amicable numbers.

Although several million examples of amicable numbers are known, it is not known whether there are infinitely many. In all known examples, the pairs of amicable numbers have the same parity and are not coprime, but proving these facts is an open problem.



## 3 Modular arithmetic

### 3.1 Congruences

**Definition 3.1.1.** Given three integers  $a, b$  and  $m$ , with  $m \geq 2$ , we say that  $a$  is congruent to  $b$  modulo  $m$ , denoted by

$$a \equiv b \pmod{m},$$

if  $a - b$  is divisible by  $m$  (or, in other words, if  $a$  and  $b$  leave the same remainder when divided by  $m$ ).

The following reformulation of congruence is often useful: we have  $a \equiv b \pmod{m}$  if and only if there exists an integer  $k$  such that  $a = b + mk$ .

**Remark 3.1.2.** Given an integer  $a$ , we may divide  $a$  by  $m$  with remainder, to write  $a = mk + r$  for some  $k \in \mathbb{Z}$  and an integer  $0 \leq r < |a|$ . Then we have  $a \equiv r \pmod{m}$ , and we call  $r$  the *least residue of  $a$  modulo  $m$* .

We collect some useful rules for computing modulo  $m$ . The proofs are straightforward verifications, which we omit.

**Lemma 3.1.3.** 1. If  $a \equiv b \pmod{m}$  and  $c \equiv d \pmod{m}$ , then  $a + c \equiv b + d \pmod{m}$  and  $ac \equiv bd \pmod{m}$ .

2. If  $a \equiv b \pmod{m}$ , then  $a^n \equiv b^n \pmod{m}$  for any  $n \in \mathbb{N}$ .

3. If  $a \equiv b \pmod{m_i}$  for  $i = 1, 2, \dots, k$ , where  $m_1, \dots, m_k$  are pairwise coprime, then  $a \equiv b \pmod{m}$ , where  $m = \prod_{i=1}^k m_i$ .

In modular arithmetic, the cancellation law does not necessarily hold, that is,  $ac \equiv bc \pmod{m}$  does not imply  $a \equiv b \pmod{m}$ . However, we have the following:

**Proposition 3.1.4.** If  $ac \equiv bc \pmod{m}$  then  $a \equiv b \pmod{m/d}$ , where  $d = \gcd(c, m)$ . In particular, if  $\gcd(c, m) = 1$ , then  $ac \equiv bc \pmod{m}$  implies  $a \equiv b \pmod{m}$ .

*Proof.* If  $ac \equiv bc \pmod{m}$ , then there exists an integer  $k$  such that  $ac - bc = km$ . Let  $d = \gcd(c, m)$ . Then

$$(a - b)(c/d) = k(m/d), \quad \gcd(c/d, m/d) = 1.$$

Hence  $m/d$  divides  $a - b$ , or, equivalently  $a \equiv b \pmod{m/d}$ . □

It is easy to check that congruence modulo  $m$  defines an equivalence relation on  $\mathbb{Z}$ , and therefore partitions  $\mathbb{Z}$  into *residue classes*

$$a + m\mathbb{Z} := \{a + mk : k \in \mathbb{Z}\}.$$

### 3 Modular arithmetic

The set of all residue classes modulo  $m$  is denoted by

$$\mathbb{Z}/m\mathbb{Z} := \{a + m\mathbb{Z} : a \in \mathbb{Z}\}.$$

The following result is easy to prove.

**Proposition 3.1.5.** *A system of representatives for  $\mathbb{Z}/m\mathbb{Z}$  is given by the least residues modulo  $m$ ,  $\{0, 1, 2, \dots, m-1\}$ . In particular, we have  $|\mathbb{Z}/m\mathbb{Z}| = m$ .*

We can equip  $\mathbb{Z}/m\mathbb{Z}$  with a well-defined addition and multiplication by setting

$$(a + m\mathbb{Z}) + (b + m\mathbb{Z}) := (a + b) + m\mathbb{Z}, \quad (a + m\mathbb{Z}) \cdot (b + m\mathbb{Z}) = (ab) + m\mathbb{Z}.$$

Moreover,  $\mathbb{Z}/m\mathbb{Z}$  defines a ring with these operations, called the *residue class ring modulo  $m$* . The *unit group* of the ring  $\mathbb{Z}/m\mathbb{Z}$  is by definition given by

$$(\mathbb{Z}/m\mathbb{Z})^* := \{a + m\mathbb{Z} : \text{there exists } b \in \mathbb{Z} \text{ such that } ab \equiv 1 \pmod{m}\}.$$

It is a group under multiplication, but not closed under addition.

**Proposition 3.1.6.** *We have*

$$(\mathbb{Z}/m\mathbb{Z})^* := \{a + m\mathbb{Z} : a \in \mathbb{Z}, \gcd(a, m) = 1\}.$$

Moreover, a system of representatives for  $(\mathbb{Z}/m\mathbb{Z})^*$  is given by

$$\{1 \leq a < m : \gcd(a, m) = 1\},$$

and we have  $|(\mathbb{Z}/m\mathbb{Z})^*| = \varphi(m)$ , with Euler's totient function  $\varphi$ .

*Proof.* If there exists some  $b \in \mathbb{Z}$  such that  $ab \equiv 1 \pmod{m}$ , then we have  $ab - 1 = km$  for some  $k \in \mathbb{Z}$ . Hence  $\gcd(a, m)$  divides 1, which means  $\gcd(a, m) = 1$ . Conversely, if  $\gcd(a, m) = 1$ , then there exist  $b, k \in \mathbb{Z}$  such that  $ab + mk = 1$ , which means that  $ab \equiv 1 \pmod{m}$ .  $\square$

**Definition 3.1.7.** An element  $a \in \mathbb{Z}$  with  $\gcd(a, m) = 1$  is called *invertible modulo  $m$* , and any  $b \in \mathbb{Z}$  with  $ab \equiv 1 \pmod{m}$  is called a *(multiplicative) inverse of  $a$  modulo  $m$* .

**Remark 3.1.8.** The residue class of the inverse of  $a$  modulo  $m$  is unique, and is sometimes denoted by  $\bar{a}$  or  $a^{-1}$  (which should not be confused with the rational number  $\frac{1}{a}$ ). Note that the multiplicative inverse of  $a$  can be determined using the Euclidean algorithm by finding  $b, k \in \mathbb{Z}$  with  $ab + mk = 1$ .

Using the inverse modulo  $m$  we can solve *linear congruences* of the form

$$ax \equiv b \pmod{m}.$$

Indeed, if  $\gcd(a, m) = 1$ , then  $x \equiv a^{-1}b \pmod{m}$  solves the above linear congruence. Hence, by multiplying with  $a^{-1}$ , we can restrict our attention to linear congruences of the form  $x \equiv b \pmod{m}$ . In order to solve *systems of linear congruences*, we can use the Chinese Remainder Theorem.

**Theorem 3.1.9** (Chinese Remainder Theorem). *Let  $m_1, \dots, m_k$  be pairwise coprime moduli, and let  $a_1, \dots, a_k$  be integers. Then the system of linear congruences*

$$\begin{aligned} x &\equiv a_1 \pmod{m_1} \\ x &\equiv a_2 \pmod{m_2} \\ &\vdots \\ x &\equiv a_k \pmod{m_k} \end{aligned}$$

has a unique solution modulo  $m = \prod_{i=1}^k m_i$ .

*Proof.* We first show that a solution exists. For  $i = 1, 2, \dots, k$  let

$$M_i = \frac{m}{m_i}.$$

Since  $m_1, \dots, m_k$  are pairwise coprime, we have  $\gcd(M_i, m_i) = 1$ . Hence,  $M_i$  is invertible modulo  $m_i$ , so we can choose an integer  $\overline{M}_i$  such that

$$M_i \overline{M}_i \equiv 1 \pmod{m_i},$$

that is,  $\overline{M}_i$  is an inverse of  $M_i$  modulo  $m_i$ . Now we put

$$x = \sum_{i=1}^k M_i \overline{M}_i a_i.$$

Since  $M_i$  is divisible by all  $m_j$  with  $i \neq j$ , computing modulo  $m_i$  we find

$$x \equiv M_i \overline{M}_i a_i \equiv a_i \pmod{m_i},$$

as desired.

Now we show that  $x$  is unique modulo  $m$ . Suppose  $x'$  is another solution of the system. Then we have

$$x \equiv x' \equiv a_i \pmod{m_i} \quad \text{for all } i \in \{1, \dots, k\}.$$

Hence  $m_i \mid (x - x')$  for all  $i$ . Since the  $m_i$  are pairwise coprime, this implies that  $m \mid (x - x')$ , so  $x \equiv x' \pmod{m}$ . Thus  $x$  is unique modulo  $m$ .  $\square$

**Remark 3.1.10.** The Chinese Remainder Theorem may be stated in more algebraic terms as follows: if  $m_1, \dots, m_k$  are pairwise coprime positive integers, and  $m = \prod_{i=1}^k m_i$ , then we have a ring isomorphism

$$\mathbb{Z}/m\mathbb{Z} \cong \mathbb{Z}/m_1\mathbb{Z} \times \cdots \times \mathbb{Z}/m_k\mathbb{Z}.$$

## 3.2 Fermat's Little Theorem

**Theorem 3.2.1** (Euler-Fermat Theorem). *If  $\gcd(a, m) = 1$ , then*

$$a^{\varphi(m)} \equiv 1 \pmod{m}.$$

### 3 Modular arithmetic

*Proof.* Let  $\{a_1, a_2, \dots, a_{\varphi(m)}\}$  be a system of representatives for  $(\mathbb{Z}/m\mathbb{Z})^*$ . Since  $\gcd(a, m) = 1$ , the multiplication-by- $a$  map  $x \mapsto ax$  is a bijection of  $(\mathbb{Z}/m\mathbb{Z})^*$ . Hence, for each  $i \in \{1, \dots, \varphi(m)\}$  there exists some  $j \in \{1, \dots, \varphi(m)\}$  such that

$$a \cdot a_i \equiv a_j \pmod{m}.$$

Thus we have

$$\prod_{i=1}^{\varphi(m)} a \cdot a_i \equiv \prod_{j=1}^{\varphi(m)} a_j \pmod{m},$$

or

$$a^{\varphi(m)} \prod_{i=1}^{\varphi(m)} a_i \equiv \prod_{j=1}^{\varphi(m)} a_j \pmod{m}$$

Since  $\gcd(a_i, m) = 1$  for all  $i$ , we can cancel  $\prod_{i=1}^{\varphi(m)} a_i$  on both sides, to find  $a^{\varphi(m)} \equiv 1 \pmod{m}$ .  $\square$

**Corollary 3.2.2** (Fermat's Little Theorem). *If  $p$  is a prime, and  $\gcd(a, p) = 1$ , then*

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

**Remark 3.2.3.** 1. The Euler-Fermat theorem can be used to compute the inverse of  $a \pmod{m}$  if  $\gcd(a, m) = 1$ : Since  $a^{\varphi(m)} \equiv 1 \pmod{m}$  and  $\varphi(m) \geq 2$  (since we assume that  $m \geq 2$ ), we have

$$a^{-1} \equiv a^{\varphi(m)-1} \pmod{m},$$

so computing the inverse  $a^{-1}$  is the same as computing the power  $a^{\varphi(m)-1}$  modulo  $m$ .

For example, let us compute the inverse of  $a = 7 \pmod{12}$ . In the exercise we proved the formula

$$\varphi(m) = m \prod_{p|m} \left(1 - \frac{1}{p}\right),$$

so we have  $\varphi(12) = 12 \cdot (1 - \frac{1}{2})(1 - \frac{1}{3}) = 4$ . Hence we obtain that  $a^{\varphi(m)-1} \equiv 7^3 \equiv 7 \pmod{12}$  is the inverse of 7 modulo 12. Indeed,  $7 \cdot 7 = 49 \equiv 1 \pmod{12}$ .

2. The Euler-Fermat theorem is helpful to compute large powers modulo  $m$ . For example, let us compute  $5^{215} \pmod{7}$ . Since  $\varphi(7) = 6$ , we have  $5^6 \equiv 1 \pmod{7}$ . Hence, if we divide 215 by 6 with remainder,  $215 = 35 \cdot 6 + 5$ , we get

$$5^{215} \equiv \underbrace{(5^6)^{35}}_{\equiv 1} \cdot 5^5 \equiv 5^5 \equiv (25)^2 \cdot 5 \equiv 4^2 \cdot 5 \equiv 2 \cdot 5 \equiv 10 \equiv 3 \pmod{7}.$$

Note that 35 did not play any role. Indeed, the above argument shows that we have

$$a^n \equiv a^{n \pmod{\varphi(m)}} \pmod{m}$$

if  $\gcd(a, m) = 1$ .



### 3.3 The theorems of Lagrange, Wilson, and Wolstenholme

In order to study polynomial congruences modulo primes  $p$ , the following general result about roots of polynomials modulo  $p$  is very useful.

**Theorem 3.3.1** (Lagrange's Theorem). *Let  $p$  be a prime, and let  $f(x) \in \mathbb{Z}[x]$  be a polynomial with integer coefficients such that not all coefficients of  $f$  are divisible by  $p$ . Then  $f(x)$  has at most  $\deg(f)$  incongruent roots modulo  $p$ .*

*Proof.* We prove the theorem by induction on  $n = \deg(f) \geq 1$ . For  $n = 1$  the equation

$$ax + b \equiv 0 \pmod{p}$$

is equivalent to  $ax \equiv -b \pmod{p}$ . If  $p \mid a$ , then  $p$  does not divide  $b$ , so  $ax \equiv -b \pmod{p}$  has no solutions. If  $p \nmid a$ , then  $a$  is invertible modulo  $p$ , so we have the unique solution  $x \equiv -a^{-1}b \pmod{p}$ , where  $a^{-1}$  is any integer with  $aa^{-1} \equiv 1 \pmod{p}$ .

Now fix some  $n \geq 1$  and assume that the theorem holds for all polynomials of degree  $\leq n$ . Let

$$f(x) = a_{n+1}x^{n+1} + \cdots + a_0 \in \mathbb{Z}[x]$$

be a polynomial of degree  $n + 1$  such that not all coefficients of  $f$  are divisible by  $p$ . If  $f(x)$  does not have any roots modulo  $p$ , there is nothing to show, so we can assume that  $f(r) \equiv 0 \pmod{p}$  for some  $r \in \mathbb{Z}$ . Since

$$x^{m+1} - r^{m+1} = (x - r) \sum_{k=0}^m x^k r^{m-k}$$

for any  $m \geq 1$ , we can write

$$\begin{aligned} f(x) &\equiv f(x) - f(r) \\ &\equiv a_{n+1}(x^{n+1} - r^{n+1}) + a_n(x^n - r^n) + \cdots + a_1(x - r) \\ &\equiv g(x)(x - r) \pmod{p}, \end{aligned}$$

where  $g(x) \in \mathbb{Z}[x]$  has degree  $\leq n$  and not all of its coefficients are divisible by  $p$  (otherwise the coefficients of  $f$  would all be divisible by  $p$ ). By induction hypothesis,  $g(x)$  has at most  $n$  roots modulo  $p$ , and  $x - r$  has precisely 1 root, so  $f(x)$  has at most  $n + 1$  roots modulo  $p$ .  $\square$

In the next section we will use Lagrange's Theorem to study quadratic congruences. As two standard applications we prove Wilson's Theorem and Wolstenholme's Theorem.

**Theorem 3.3.2** (Wilson's Theorem). *A natural number  $n > 1$  is prime if and only if*

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

*Proof.* Suppose that  $p$  is prime. We consider the two polynomials

$$g(x) = x^{p-1} - 1, \quad h(x) = (x - 1)(x - 2) \cdots (x - (p - 1))$$

with integer coefficients and degree  $p - 1$ . By Fermat's Little Theorem, we have  $g(x) \equiv 0 \pmod{p}$  for  $x = 1, 2, \dots, p - 1$ , and by construction we have  $h(x) \equiv 0 \pmod{p}$  for  $x = 1, 2, \dots, p - 1$ , as well. Now

$$f(x) = g(x) - h(x)$$

### 3 Modular arithmetic

is a polynomial of degree at most  $p - 2$  (since  $g(x)$  and  $h(x)$  have the same leading term), and the equation  $f(x) \equiv 0 \pmod{p}$  has  $p - 1$  incongruent solutions. By Lagrange's Theorem, this means that every coefficient of  $f(x)$  is divisible by  $p$ . In particular, the constant term of  $f(x)$  is divisible by  $p$ , which means that

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

If  $p$  is odd, then  $(-1)^{p-1} = 1$ , and if  $p = 2$ , then  $(-1)^{p-1} = -1 \equiv 1 \pmod{2}$ , so we obtain  $-1 \equiv (p-1)! \pmod{p}$  for any prime  $p$ , as claimed.

Conversely, if  $n$  is composite and  $n > 4$ , then  $(n-1)! \equiv 0 \pmod{n}$ . We leave this as an exercise to the reader. This implies that  $(n-1)! \not\equiv -1 \pmod{n}$ . For  $n = 4$  we have  $(n-1)! = 6 \equiv 2 \pmod{4}$ . This finishes the proof.  $\square$

A prime  $p$  is called a *Wilson prime* if it satisfies the stronger congruence

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

The only known Wilson primes are 5, 13, and 563, but it is conjectured that there are infinitely many. However, the next one must be bigger than  $10^{13}$ .

**Theorem 3.3.3** (Wolstenholme's Theorem). *If  $p > 3$  is a prime, then the numerator of the rational number*

$$1 + \frac{1}{2} + \frac{1}{3} + \cdots + \frac{1}{p-1}$$

*is divisible by  $p^2$ , and the numerator of the rational number*

$$1 + \frac{1}{2^2} + \frac{1}{3^2} + \cdots + \frac{1}{(p-1)^2}.$$

*is divisible by  $p$ . Moreover, we have the congruence*

$$\binom{2p-1}{p-1} \equiv 1 \pmod{p^3}.$$

*for every prime  $p > 3$ .*

**Example 3.3.4.** For  $p = 5$  we have  $1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} = \frac{25}{12}$  and  $1 + \frac{1}{4} + \frac{1}{9} + \frac{1}{16} = \frac{205}{144}$ , and  $\binom{2 \cdot 5 - 1}{5 - 1} = 126 \equiv 1 \pmod{5^3}$ .

*Proof.* We have seen in the proof of Wilson's Theorem that the coefficients  $a_1, \dots, a_{p-2}$  appearing in the polynomial

$$h(x) = (x-1)(x-2)\cdots(x-(p-1)) = x^{p-1} + a_{p-2}x^{p-2} + \cdots + a_1x + (p-1)!$$

are all divisible by  $p$ . From the product for  $h(x)$  we see that

$$a_1 = -(p-1)! \left( 1 + \frac{1}{2} + \frac{1}{3} + \cdots + \frac{1}{p-1} \right).$$

Since  $(p-1)!$  is not divisible by  $p$ , it suffices to show that  $a_1$  is divisible by  $p^2$ . To this end, note that  $h(p) = (p-1)!$ , so we find

$$(p-1)! = p^{p-1} + a_{p-2}p^{p-2} + \cdots + a_2p^2 + a_1p + (p-1)!$$

### 3.3 The theorems of Lagrange, Wilson, and Wolstenholme

Cancelling  $(p-1)!$  on both sides, and reducing modulo  $p^3$  (here we use that  $p > 3$ , so  $p^{p-1} \equiv 0 \pmod{p^3}$ ), we find

$$0 \equiv a_1 p \pmod{p^3},$$

which shows that  $a_1 \equiv 0 \pmod{p^2}$ . This finishes the proof of the first claim.

For the second claim, we write

$$1 + \frac{1}{2^2} + \frac{1}{3^2} + \cdots + \frac{1}{(p-1)^2} = \frac{1}{(p-1)!^2} A,$$

with the integer

$$A = \frac{(p-1)!^2}{1^2} + \frac{(p-1)!^2}{2^2} + \frac{(p-1)!^2}{3^2} + \cdots + \frac{(p-1)!^2}{(p-1)^2}.$$

Since  $(p-1)!^2$  is coprime to  $p$ , it suffices to show that  $A$  is divisible by  $p$ . We claim that

$$\frac{(p-1)!^2}{a^2} \equiv (a^{-1})^2 \pmod{p}$$

for  $\gcd(a, p) = 1$ . Indeed, we have

$$a^2 \cdot \frac{(p-1)!^2}{a^2} \equiv (p-1)!^2 \equiv (-1)^2 \equiv 1 \pmod{p}$$

by Willson's Theorem. Hence, we can rewrite

$$A \equiv (1^{-1})^2 + (2^{-1})^2 + (3^{-1})^2 + \cdots + ((p-1)^{-1})^2 \pmod{p}.$$

Since the map  $a \mapsto a^{-1}$  just permutes the elements of  $(\mathbb{Z}/p\mathbb{Z})^*$ , we obtain

$$A \equiv 1^2 + 2^2 + 3^2 + \cdots + (p-1)^2 \pmod{p}.$$

Now we have the formula  $1^2 + 2^2 + 3^2 + \cdots + (p-1)^2 = \frac{(p-1)p(2p-1)}{6}$ , so

$$A \equiv \frac{(p-1)p(2p-1)}{6} \pmod{p}.$$

Hence, for  $p > 3$ , we see that  $A$  is divisible by  $p$ . This finishes the proof of the second claim. The third claim easily follows from the first one and will be treated in the exercises.  $\square$

The following “converse” of Wolstenholme's Theorem is conjectured to be true:

**Conjecture 3.3.5** (Wolstenholme). *If  $\binom{2n-1}{n-1} \equiv 1 \pmod{n^3}$  for a natural number  $n$ , then  $n$  must be prime.*

The conjecture has been verified numerically for  $n \leq 10^9$ .

A prime  $p > 3$  is called a *Wolstenholme prime* if it satisfies the stronger congruence

$$\binom{2p-1}{p-1} \equiv 1 \pmod{p^4}.$$

So far, the only known Wolstenholme primes are 16843 (found in 1964) and 2124679 (found in 1993), and these are the only ones up to  $10^9$ . However, it is conjectured that infinitely many Wolstenholme primes exist.

### 3.4 Quadratic residues

In this section, we study quadratic congruences of the form  $x^2 \equiv a \pmod{p}$  for odd primes  $p$ .

**Definition 3.4.1.** Let  $p$  be a prime, and let  $a$  be an integer with  $\gcd(a, p) = 1$ . We say that  $a$  is a *quadratic residue modulo  $p$*  if the *quadratic congruence*

$$x^2 \equiv a \pmod{p},$$

has a solution. Otherwise,  $a$  is called a *quadratic nonresidue*.

If  $p$  is an odd prime and  $\gcd(a, p) = 1$ , we define the *Legendre symbol*

$$\left(\frac{a}{p}\right) = \begin{cases} 1 & \text{if } a \text{ is a quadratic residue modulo } p, \\ -1 & \text{if } a \text{ is a quadratic nonresidue modulo } p. \end{cases}$$

**Remark 3.4.2.** If  $a \equiv 0 \pmod{p}$ , then  $x^2 \equiv a \pmod{p}$  is trivially solvable for every  $p$ . In this case, the Legendre symbol  $\left(\frac{a}{p}\right)$  is usually defined to be 0. We exclude this case for simplicity, and do not call 0 (or any multiple of  $p$ ) a quadratic residue modulo  $p$  in the following.

Note that, if  $p$  is an odd prime, and  $a, b$  are coprime to  $p$  with  $a \equiv b \pmod{p}$ , then

$$\left(\frac{a}{p}\right) = \left(\frac{b}{p}\right).$$

We first show that exactly half of the elements in  $(\mathbb{Z}/p\mathbb{Z})^*$  are quadratic residues, and the other half are quadratic non-residues.

**Theorem 3.4.3.** *If  $p$  is an odd prime, then there are precisely  $(p-1)/2$  incongruent quadratic residues modulo  $p$ , given by*

$$1^2, 2^2, \dots, \left(\frac{p-1}{2}\right)^2 \pmod{p}.$$

*Proof.* The quadratic residues modulo  $p$  are given by the incongruent values of  $x^2 \pmod{p}$  for  $x \in \{1, \dots, p-1\}$ . Since

$$x^2 \equiv (p-x)^2 \pmod{p},$$

the squares of the numbers in the sets

$$\{1, 2, \dots, (p-1)/2\}, \quad \{(p-1)/2 + 1, \dots, p-1\}$$

are congruent in pairs. So we only need to consider the numbers in the first set. We claim that their squares  $1^2, 2^2, \dots, \left(\frac{p-1}{2}\right)^2$  are all incongruent modulo  $p$ . Indeed, otherwise the equation  $x^2 \equiv a \pmod{p}$  would have four incongruent solutions mod  $p$  for some  $a$ , contradicting Lagrange's Theorem. Thus, the quadratic residues modulo  $p$  are precisely the integers  $1^2, 2^2, \dots, \left(\frac{p-1}{2}\right)^2$ .  $\square$

Next, we give a refinement of Fermat's Little Theorem, using the Legendre symbol. Recall that Fermat's Little Theorem states that, for  $p$  an odd prime and  $\gcd(a, p) = 1$ , we have

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

Now from

$$a^{p-1} - 1 = (a^{(p-1)/2} - 1)(a^{(p-1)/2} + 1) \equiv 0 \pmod{p}$$

it follows that either

$$a^{(p-1)/2} \equiv +1 \pmod{p} \quad \text{or} \quad a^{(p-1)/2} \equiv -1 \pmod{p}.$$

The sign on the right-hand side is determined by the Legendre symbol:

**Theorem 3.4.4** (Euler's criterion). *If  $p$  is an odd prime and  $\gcd(a, p) = 1$ , then*

$$a^{(p-1)/2} \equiv \left(\frac{a}{p}\right) \pmod{p}.$$

*Proof.* By Fermat's Little Theorem, the polynomial  $x^{p-1} - 1$  has precisely  $p - 1$  incongruent roots modulo  $p$ , namely those in  $(\mathbb{Z}/p\mathbb{Z})^*$ . Since  $p - 1$  is even, we can factor the polynomial as

$$x^{p-1} - 1 = (x^{(p-1)/2} - 1)(x^{(p-1)/2} + 1).$$

If  $a$  is a quadratic residue modulo  $p$ , then there is some integer  $b$  with  $b^2 \equiv a \pmod{p}$ . Hence

$$a^{(p-1)/2} \equiv b^{p-1} \equiv 1 \equiv \left(\frac{a}{p}\right) \pmod{p}$$

by Fermat's Little Theorem. This means that each quadratic residue is a zero of the polynomial  $x^{(p-1)/2} - 1$  modulo  $p$ . We have seen that precisely  $(p - 1)/2$  of the elements of  $(\mathbb{Z}/p\mathbb{Z})^*$  are quadratic residues, so the polynomial  $x^{(p-1)/2} - 1$  has precisely the quadratic residues as roots modulo  $p$  by Lagrange's Theorem. Hence, the quadratic nonresidues must be the zeros of  $x^{(p-1)/2} + 1$  modulo  $p$ . In other words, if  $a$  is a quadratic nonresidue, then

$$a^{(p-1)/2} \equiv -1 \equiv \left(\frac{a}{p}\right) \pmod{p}$$

This finishes the proof. □

Note that Euler's criterion gives a method to determine  $\left(\frac{a}{p}\right)$ : just compute  $a^{(p-1)/2} \pmod{p}$ , e.g. by repeated squaring. The next two results follow from Euler's criterion, and will be left as an exercise.

**Theorem 3.4.5.** *If  $p$  is an odd prime, then*

$$\left(\frac{-1}{p}\right) = (-1)^{(p-1)/2} = \begin{cases} 1 & \text{if } p \equiv 1 \pmod{4}, \\ -1 & \text{if } p \equiv 3 \pmod{4}. \end{cases}$$

**Theorem 3.4.6.** *If  $p$  is an odd prime and  $p$  does not divide  $ab$ , then*

$$\left(\frac{ab}{p}\right) = \left(\frac{a}{p}\right) \left(\frac{b}{p}\right).$$

*In other words, the function  $a \mapsto \left(\frac{a}{p}\right)$  is completely multiplicative.*

### 3 Modular arithmetic

Note that the above multiplicativity implies the peculiar fact that the product of two quadratic nonresidues is a quadratic residue!

**Theorem 3.4.7** (Gauss's Lemma). *Let  $p$  be an odd prime and  $\gcd(a, p) = 1$ . Consider the  $(p-1)/2$  numbers*

$$a, 2a, 3a, \dots, \frac{p-1}{2}a,$$

*and their least residues modulo  $p$ . Let  $s$  be the number of these least residues larger than  $p/2$ . Then we have*

$$\left(\frac{a}{p}\right) = (-1)^s.$$

*Proof.* Let  $S$  denote the set of least residues modulo  $p$  of the numbers  $a, 2a, \dots, \frac{p-1}{2}a$ . Note that the elements of  $S$  are pairwise distinct. Let  $s$  be the number of elements in  $S$  that exceed  $p/2$ , and put  $r = \frac{p-1}{2} - s$ . Relabel the elements of  $S$  as

$$a_1, a_2, \dots, a_r, \quad b_1, b_2, \dots, b_s,$$

where  $a_i < p/2$  and  $b_j > p/2$ . Since the elements of  $S$  are the least residues of  $a, 2a, \dots, \frac{p-1}{2}a$ , we have

$$\left(\prod_{i=1}^r a_i\right) \left(\prod_{j=1}^s b_j\right) \equiv \prod_{k=1}^{\frac{p-1}{2}} (ka) \equiv a^{(p-1)/2} \left(\frac{p-1}{2}\right)! \pmod{p}. \quad (3.4.1)$$

Consider the set  $T$  consisting of the  $\frac{p-1}{2}$  integers

$$a_1, a_2, \dots, a_r, \quad p - b_1, p - b_2, \dots, p - b_s.$$

Then all elements of  $T$  lie between 1 and  $\frac{p-1}{2}$ . Moreover, we claim that the integers in  $T$  are pairwise distinct, and hence  $T$  contains precisely the numbers  $1, 2, \dots, \frac{p-1}{2}$ . Since the  $a_i$  and  $b_j$  are pairwise distinct, it suffices to show that we cannot have  $a_i \equiv p - b_j \pmod{p}$  for any  $i, j$ . Indeed, assume  $a_i \equiv p - b_j \pmod{p}$ , and write  $a_i = ha$  and  $b_j = ka$  with  $1 \leq h, k \leq \frac{p-1}{2}$  and  $h \neq k$ . Then

$$0 \equiv p \equiv a_i + b_j \equiv (h+k)a \pmod{p}.$$

Since  $p$  does not divide  $a$ , this means that  $p$  must divide  $h+k$ . But we have  $0 < h+k < p$ , so this is impossible. This shows that  $T$  consists of the numbers  $1, 2, \dots, \frac{p-1}{2}$ .

Taking into account (3.4.1) we obtain

$$\begin{aligned} \left(\frac{p-1}{2}\right)! &\equiv \prod_{n \in T} n \equiv \left(\prod_{i=1}^r a_i\right) \left(\prod_{j=1}^s (p - b_j)\right) \equiv (-1)^s \left(\prod_{i=1}^r a_i\right) \left(\prod_{j=1}^s b_j\right) \\ &\equiv (-1)^s a^{(p-1)/2} \left(\frac{p-1}{2}\right)! \pmod{p} \end{aligned}$$

Cancelling  $\left(\frac{p-1}{2}\right)!$  from both sides yields

$$1 \equiv (-1)^s a^{(p-1)/2} \pmod{p}.$$

From Euler's criterion, we obtain  $\left(\frac{a}{p}\right) \equiv (-1)^s \pmod{p}$ , as claimed.  $\square$

**Example 3.4.8.** Let  $p = 7$  and  $a = 3$ . In order to apply Gauss's Lemma, we consider the first  $(p - 1)/2 = 3$  multiples of 3,

$$3, 6, 9,$$

and their least positive residues modulo 7,

$$3, 6, 2,$$

of which only one is larger than  $p/2 = 3.5$ . Hence Gauss's Lemma tells us that  $\left(\frac{3}{7}\right) = -1$ . Indeed, the squares modulo 7 are 1, 2, 4, so 3 is a quadratic nonresidue modulo 7.

The next result follows from Gauss's Lemma, and will be treated in the exercises.

**Theorem 3.4.9.** *If  $p$  is an odd prime, then*

$$\left(\frac{2}{p}\right) = (-1)^{(p^2-1)/8} = \begin{cases} 1 & \text{if } p \equiv \pm 1 \pmod{8}, \\ -1 & \text{if } p \equiv \pm 3 \pmod{8}. \end{cases}$$

The main result about quadratic congruences and the Legendre symbol is the *quadratic reciprocity law*.

**Theorem 3.4.10** (Gauss's quadratic reciprocity law). *If  $p$  and  $q$  are distinct odd primes, then*

$$\left(\frac{p}{q}\right) \left(\frac{q}{p}\right) = (-1)^{\frac{p-1}{2} \frac{q-1}{2}}.$$

It was first proved by Gauss in 1801. Gauss found (at least) 8 different proofs, and until today people have found more than 240 proofs of the quadratic reciprocity law. We will not give a proof here, but refer to "Proofs from THE BOOK" for an elementary proof.

The two statements

1.  $\left(\frac{-1}{p}\right) = (-1)^{(p-1)/2} = \begin{cases} 1 & \text{if } p \equiv 1 \pmod{4}, \\ -1 & \text{if } p \equiv 3 \pmod{4}, \end{cases}$
2.  $\left(\frac{2}{p}\right) = (-1)^{(p^2-1)/8} = \begin{cases} 1 & \text{if } p \equiv \pm 1 \pmod{8}, \\ -1 & \text{if } p \equiv \pm 3 \pmod{8}, \end{cases}$

are often called the *first and second supplement* to the quadratic reciprocity law.

Using the quadratic reciprocity law and its supplements one can compute the Legendre symbol.

**Example 3.4.11.** In order to determine whether the quadratic congruence  $x^2 \equiv 24 \pmod{101}$  has a solution, we compute the Legendre symbol  $\left(\frac{24}{101}\right)$ . Note that 101 is prime, and  $24 = 3 \cdot 2^3$ . We first use that the Legendre symbol is completely multiplicative in the "numerator",

$$\left(\frac{24}{101}\right) = \left(\frac{3}{101}\right) \left(\frac{2}{101}\right)^3 = \left(\frac{3}{101}\right) \left(\frac{2}{101}\right),$$

where we used that  $\left(\frac{a}{p}\right)^2 = 1$  if  $\gcd(a, p) = 1$ . Now, by quadratic reciprocity, we have

$$\left(\frac{3}{101}\right) = (-1)^{\frac{3-1}{2} \frac{101-1}{2}} \left(\frac{101}{3}\right) = \left(\frac{101}{3}\right) = \left(\frac{-1}{3}\right),$$

### 3 Modular arithmetic

where we used that we may reduce the “numerator” modulo the “denominator” in the Legendre symbol. It remains to compute  $\left(\frac{2}{101}\right)$  and  $\left(\frac{-1}{3}\right)$ , which can be done using the supplements:

$$\left(\frac{2}{101}\right) = -1, \quad \left(\frac{-1}{3}\right) = (-1)^{\frac{3-1}{2}} = -1,$$

since  $101 \equiv -3 \pmod{8}$ . In total, we find

$$\left(\frac{24}{101}\right) = \left(\frac{3}{101}\right) \left(\frac{2}{101}\right) = \left(\frac{-1}{3}\right) \left(\frac{2}{101}\right) = (-1) \cdot (-1) = 1,$$

so 24 is a square modulo 101. Indeed, by trying  $x = 1, 2, 3, \dots$  we find that  $x = 23$  solves  $x^2 \equiv 24 \pmod{101}$ .

**Remark 3.4.12.** The Legendre symbol can be generalized to cases where the denominator is composite. If  $m = p_1^{\nu_1} \cdots p_r^{\nu_r}$  is an *odd* natural number and  $\gcd(a, m) = 1$ , the *Jacobi symbol* is defined by

$$\left(\frac{a}{m}\right) = \left(\frac{a}{p_1}\right)^{\nu_1} \cdots \left(\frac{a}{p_r}\right)^{\nu_r}.$$

Note that the Jacobi symbol  $\left(\frac{a}{m}\right)$  no longer tells us whether  $a$  is a quadratic residue modulo  $m$ , that is, we may have  $\left(\frac{a}{m}\right) = 1$  even if  $x^2 \equiv a \pmod{m}$  is not solvable. However, if  $\left(\frac{a}{m}\right) = -1$ , then the equation  $x^2 \equiv a \pmod{m}$  has no solution.

We list some important properties of the Jacobi symbol, which are not hard to deduce from the properties of the Legendre symbol (recall that  $m$  is odd):

1.  $a \equiv b \pmod{m}$  implies  $\left(\frac{a}{m}\right) = \left(\frac{b}{m}\right)$ ,
2.  $\left(\frac{ab}{m}\right) = \left(\frac{a}{m}\right) \left(\frac{b}{m}\right)$ ,
3.  $\left(\frac{a}{mn}\right) = \left(\frac{a}{m}\right) \left(\frac{a}{n}\right)$  if  $n$  is odd,
4.  $\left(\frac{-1}{m}\right) = (-1)^{\frac{m-1}{2}}$ ,
5.  $\left(\frac{2}{m}\right) = (-1)^{\frac{m^2-1}{8}}$ ,
6.  $\left(\frac{n}{m}\right) \left(\frac{m}{n}\right) = (-1)^{\frac{n-1}{2} \frac{m-1}{2}}$  for  $m, n$  odd with  $\gcd(m, n) = 1$ .

An important application is the fact that these rules can be used to compute the Jacobi symbol  $\left(\frac{a}{m}\right)$  *without factoring*  $a$ . This is important since the factorization of large integers becomes very hard, but the Jacobi symbol can still be computed quickly.

**Example 3.4.13.** Let us compute  $\left(\frac{90}{101}\right)$  using the Jacobi symbol, without completely factoring the numerator 90. We need to pull out all factors of 2 from 90 (this is easy to do also for large numbers), and write

$$\left(\frac{90}{101}\right) = \left(\frac{2 \cdot 45}{101}\right) = \left(\frac{2}{101}\right) \left(\frac{45}{101}\right).$$



The first factor equals  $\left(\frac{2}{101}\right) = -1$  since  $101 \equiv -3 \pmod{8}$ . In order to compute the second factor, we use quadratic reciprocity,

$$\left(\frac{45}{101}\right) = (-1)^{\frac{44}{2} \cdot \frac{100}{2}} \left(\frac{101}{45}\right) = \left(\frac{11}{45}\right),$$

where we reduced the numerator modulo 45 in the second step. We can again apply quadratic reciprocity,

$$\left(\frac{11}{45}\right) = (-1)^{\frac{10}{2} \cdot \frac{44}{2}} \left(\frac{45}{11}\right) = \left(\frac{1}{11}\right) = 1.$$

In total, we find  $\left(\frac{90}{101}\right) = -1$ .

## 3.5 Applications to cryptography

### 3.5.1 Primality tests

In cryptographic applications it is important to be able to efficiently generate very large primes. The general idea is to take a large random number  $n$  and then “test” whether  $n$  is prime. For the sake of efficiency, such a test does typically not verify that a number is prime for certain, but only that it is “very probably” prime. Here we describe two simple tests which are based on Fermat’s Little Theorem and Euler’s criterion.

Fermat’s Little Theorem states that, for a prime  $p$  and  $a \in \mathbb{Z}$  with  $\gcd(a, p) = 1$ , we have

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

In general, this statement becomes false if we replace  $p$  with a composite number  $n$ : for “most” composite natural numbers  $n$ , there exists some  $a \in \mathbb{Z}$  with  $\gcd(a, n) = 1$  but  $a^{n-1} \not\equiv 1 \pmod{n}$ . For example, for  $n = 4$  and  $a = 2$  we have  $2^3 \equiv 8 \equiv 0 \pmod{4}$ , so  $n = 4$  cannot be prime. This fact is used in *Fermat’s primality test*.

**Algorithm 3.5.1** (Fermat’s primality test). *Given an integer  $n$ , the following algorithm either returns “ $n$  is composite” or “ $n$  is possibly prime”.*

1. Choose a random  $a \in \{1, \dots, n-1\}$  with  $\gcd(a, n) = 1$  (check this using the Euclidean algorithm).
2. If  $a^{n-1} \not\equiv 1 \pmod{n}$ , then  $n$  is composite.  
Return “ $n$  is composite”.
3. If  $a^{n-1} \equiv 1 \pmod{n}$ , we cannot be sure whether  $n$  is prime or composite.  
Return “ $n$  is possibly prime”.

It is important to note that if the test returns “ $n$  is composite”, then we can be sure that  $n$  is composite, but the test does *not* yield a factorization of  $n$  or even a non-trivial divisor.

If the algorithm returns “ $n$  is possibly prime”, we may repeat the test several times with different choices of  $a$  to become more confident that  $n$  is prime, but we can never be completely sure: there are some exceptional composite numbers that the Fermat test cannot detect.

**Definition 3.5.2.** A composite number  $n$  for which  $a^{n-1} \equiv 1 \pmod{n}$  for every  $a \in \mathbb{Z}$  with  $\gcd(a, n) = 1$  is called a *Carmichael number* (or a *pseudo-prime*).

### 3 Modular arithmetic

On a Carmichael number  $n$ , the Fermat test will always return “ $n$  is possibly prime”. The smallest Carmichael number is 561. It was shown in 1994 by Alford, Granville, and Pomerance [?] that there are infinitely many Carmichael numbers. Although they are rather rare, it was shown in 2021 by Daniel Larson that the Carmichael numbers satisfy an analogue of Bertrand’s postulate, that is, it is possible to control the gaps between two consecutive Carmichael numbers.

**Theorem 3.5.3** (Korselt’s criterion). *A composite number  $n$  is a Carmichael number if and only if  $n$  is square-free and  $(p - 1) \mid (n - 1)$  for every  $p \mid n$ .*

*Proof.* Suppose that  $n$  is square-free and  $(p - 1) \mid (n - 1)$  for every  $p \mid n$ . Let  $n = p_1 \cdots p_r$  be the prime factorization of  $n$ , with distinct primes  $p_j$ . For  $a \in \mathbb{Z}$  with  $\gcd(a, n) = 1$  we also have  $\gcd(a, p_j) = 1$ , so by Fermat’s Little Theorem

$$a^{p_j-1} \equiv 1 \pmod{p_j}.$$

From  $(p_j - 1) \mid (n - 1)$  it follows that

$$a^{n-1} \equiv 1 \pmod{p_j},$$

and the Chinese Remainder Theorem implies

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

and since  $a$  was arbitrary with  $\gcd(a, n) = 1$ , we see that  $n$  is a Carmichael number.

For the converse direction we will use some simple facts from elementary group theory. The *order*  $\text{ord}(a)$  of an element  $a \in (\mathbb{Z}/p\mathbb{Z})^*$  is the smallest positive number  $m$  such that  $a^m \equiv 1 \pmod{p}$ . The order only depends on  $a$  modulo  $p$ , and if  $a^\ell \equiv 1 \pmod{p}$  for some  $\ell \in \mathbb{Z}$ , then  $\text{ord}(a) \mid \ell$ . The group  $(\mathbb{Z}/p\mathbb{Z})^*$  is *cyclic*, which means that there exists some  $a \in (\mathbb{Z}/p\mathbb{Z})^*$ , a *generator*, such that every other element in  $(\mathbb{Z}/p\mathbb{Z})^*$  can be obtained as a suitable power of  $a$ . Each such generator has the maximal possible order  $\text{ord}(a) = p - 1$ .

Let  $n$  be a Carmichael number and write  $n = p^k n'$  with an odd prime  $p$  and some natural number  $n'$  coprime to  $p$ . Take a number  $a \in \mathbb{Z}$  with  $\gcd(a, p) = 1$  and order  $p - 1$  in  $(\mathbb{Z}/p\mathbb{Z})^*$ . In order to use that  $n$  is a Carmichael number, we would need that  $\gcd(a, n) = 1$ . By the Chinese Remainder Theorem, there exists an  $a' \in \mathbb{Z}$  such that  $a' \equiv a \pmod{p}$  and  $a' \equiv 1 \pmod{n'}$ , and since the order only depends on  $a$  modulo  $p$ ,  $a'$  still has order  $p - 1$ . Note that now we have  $\gcd(a', n) = 1$ , and since  $n$  is a Carmichael number, we have

$$a'^{n-1} \equiv 1 \pmod{n}.$$

In particular, reducing further modulo  $p$ , we have

$$a'^{n-1} \equiv 1 \pmod{p}.$$

Since  $a'$  has order  $p - 1$  in  $(\mathbb{Z}/p\mathbb{Z})^*$ , this implies that  $(p - 1) \mid (n - 1)$ .

We now show that a Carmichael number  $n$  has to be square-free. Let us assume that a prime factor  $p$  of  $n$  appears more than once, and write  $n = p^k n'$  with  $k \geq 2$  and  $\gcd(p, n') = 1$ . By the Chinese Remainder Theorem, there exists an  $a \in \mathbb{Z}$  with

$$a \equiv 1 + p \pmod{p^2}, \quad a \equiv 1 \pmod{n'}.$$

Since  $\gcd(a, n) = 1$ , and  $n$  is Carmichael number, we have

$$a^{n-1} \equiv 1 \pmod{n}.$$

Reducing modulo  $p^2$ , and using that  $a \equiv 1 + p \pmod{p^2}$ , we obtain

$$(1 + p)^{n-1} \equiv 1 \pmod{p^2}.$$

By the Binomial Theorem, we see that

$$(1 + p)^{n-1} \equiv 1 + (n-1)p \pmod{p^2},$$

and since  $n$  is divisible by  $p$ , we finally obtain

$$1 - p \equiv 1 \pmod{p^2},$$

which is a contradiction. This shows that a Carmichael number  $n$  must be square-free, and finishes the proof.  $\square$

A further restriction for Carmichael numbers is given in the following simple lemma, which follows from Korselt's criterion and will be treated in the exercises.

**Lemma 3.5.4.** *A Carmichael number  $n$  must be odd, have at least three distinct prime factors, and for primes  $p, q$  dividing  $n$  we have  $p \not\equiv 1 \pmod{q}$ .*

The Fermat primality test has the flaw that it always mistakes Carmichael numbers for primes. In order to refine the test, we can use Euler's criterion, which states that, for an odd prime  $p$  and  $\gcd(a, p) = 1$  we have

$$a^{\frac{p-1}{2}} \equiv \left(\frac{a}{p}\right) \pmod{p},$$

which leads to the following primality test.

**Algorithm 3.5.5** (Solovay-Strassen primality test). *Given an odd integer  $n$ , the following algorithm either returns “ $n$  is composite” or “ $n$  is possibly prime”.*

1. Choose a random  $a \in \{1, \dots, n-1\}$  with  $\gcd(a, n) = 1$ .
2. If  $a^{\frac{n-1}{2}} \not\equiv \left(\frac{a}{n}\right) \pmod{n}$ , then  $n$  is composite.  
Return “ $n$  is composite”.
3. If  $a^{\frac{n-1}{2}} \equiv \left(\frac{a}{n}\right) \pmod{n}$ , we cannot be sure whether  $n$  is prime or composite.  
Return “ $n$  is possibly prime”.

Note that we can assume that  $n$  is odd since even numbers are not prime and are easily identified from their last digit. Moreover, the Jacobi symbols  $\left(\frac{a}{n}\right)$  can be efficiently computed (without factoring  $n$ ) using quadratic reciprocity as discussed in the last section.

There are no “strong” pseudo-primes, i.e. no analogues of the Carmichael numbers that can fool the Solovay-Strassen test:

**Theorem 3.5.6** (Lehmer). *If  $a^{\frac{n-1}{2}} \equiv \left(\frac{a}{n}\right) \pmod{n}$  for all  $a$  with  $\gcd(a, n) = 1$ , then  $n$  must be prime.*

### 3 Modular arithmetic

*Proof.* Suppose that  $n$  is composite. Any composite number with the above property is also a Carmichael number, so by Korselt's criterion we know that  $n = p_1 \cdots p_r$  with at least two distinct odd primes  $p_1, \dots, p_r$ .

Choose some quadratic nonresidue  $b$  modulo  $p_1$ , that is,  $\left(\frac{b}{p_1}\right) = -1$ . By the Chinese Remainder Theorem, there is some  $a \in \mathbb{Z}$  with

$$a \equiv b \pmod{p_1} \quad \text{and} \quad a \equiv 1 \pmod{p_2}, \dots, a \equiv 1 \pmod{p_r}.$$

In particular, we have

$$\left(\frac{a}{n}\right) = \left(\frac{a}{p_1}\right) \left(\frac{a}{p_2}\right) \cdots \left(\frac{a}{p_r}\right) = \left(\frac{b}{p_1}\right) \left(\frac{1}{p_2}\right) \cdots \left(\frac{1}{p_r}\right) = -1,$$

and by the assumption on  $n$ , we get

$$a^{\frac{n-1}{2}} \equiv \left(\frac{a}{n}\right) \equiv -1 \pmod{n}.$$

However, reducing this modulo  $p_2$  and using  $a \equiv 1 \pmod{p_2}$  gives the contradiction

$$1 \equiv -1 \pmod{p_2}.$$

Hence  $n$  cannot be composite. □

**Remark 3.5.7.** For composite odd  $n$ , the congruence  $a^{\frac{n-1}{2}} \equiv \left(\frac{a}{n}\right) \pmod{n}$  can be satisfied for at most half of the elements  $a \in (\mathbb{Z}/n\mathbb{Z})^*$ . Using a little group theory, this easily follows from the fact that the set of all these “bad”  $a$  forms a proper subgroup of  $(\mathbb{Z}/n\mathbb{Z})^*$ , and hence the order of this subgroup has to divide the order of  $(\mathbb{Z}/n\mathbb{Z})^*$  by a theorem of Lagrange, so it can be at most half as large.

However, in practice, one never checks all possible values of  $a$ , since the possibility that  $n$  is prime if the algorithm returns “ $n$  is possibly prime” for several random choices of  $a$  goes to 1 quite fast.

There are other, more refined primality tests, e.g. the Miller-Rabin test, which rely on similar ideas from elementary number theory and are more relevant in practical cryptographic applications.

#### 3.5.2 The RSA cryptosystem

The main goal of a cryptosystem is the secure transmission of secret messages. Classical cryptosystems, such as the Caesar cipher, are *symmetric*, which means that the sender and the receiver use *the same secret key* to encrypt and decrypt messages. The drawback of such a system is that the common secret key has to be exchanged first, which is a serious risk.

Around 1970, Diffie and Hellman, and Rivest, Shamir and Adleman developed the idea of *public-key* cryptosystems, where the sender and the receiver use *different keys*, so that there is no need for a secret key exchange anymore. The sender uses the *public key* of the receiver to encrypt messages. This public key is freely available for anyone, and can only be used for encryption, but not for decryption. The receiver uses his *private key* to decrypt messages. The private key is not distributed, but remains the receiver's secret.

One of the first and still most important public-key systems is the RSA cryptosystem. Each participant in the RSA cryptosystem needs to generate his own public and private keys:

**Algorithm 3.5.8** (RSA cryptosystem - key generation). *Alice wants to send a secret message to Bob. In order to encrypt messages, she needs Bob's public key. Here is how Bob generates his public and private keys:*

1. Choose two very large random primes  $p$  and  $q$ .
2. Compute the RSA modulus  $N = pq$ .
3. Compute Euler's totient function  $\varphi(N) = (p - 1)(q - 1)$ .
4. Choose some number  $e$  with  $1 < e < \varphi(N)$  which is coprime to  $\varphi(N)$ .
5. Compute the inverse  $d$  of  $e$  modulo  $\varphi(N)$ , that is,  $1 < d < \varphi(N)$  and  $ed \equiv 1 \pmod{\varphi(N)}$ .

The public key is the pair  $(e, N)$ .

The private key is the pair  $(d, N)$ .

The values  $p, q$ , and  $\varphi(N)$  will not be used anymore, and should be destroyed, since the knowledge of  $p, q$  or  $\varphi(N)$  would enable an attacker to compute the private key and decrypt secret messages. Note that the computation of the inverse  $d$  of  $e$  modulo  $\varphi(N)$  can be quickly done using the Euclidean algorithm.

**Example 3.5.9.** Let us generate some public and private RSA keys. We pick the primes<sup>1</sup>

$$p = 823, \quad q = 409,$$

and form the RSA modulus  $N = pq$  and its  $\varphi$ -function,

$$N = pq = 823 \cdot 409 = 336607, \quad \varphi(N) = (p - 1)(q - 1) = 822 \cdot 408 = 335376.$$

Now we choose the public key. We randomly pick<sup>2</sup>

$$e = 2147,$$

which satisfies the requirements  $1 < e < \varphi(N)$  and  $\gcd(e, \varphi(N)) = 1$ . Inverting  $e$  modulo  $\varphi(N) = 335376$ , we obtain the private key

$$d = 280235.$$

The public key  $e = 2147$  and the RSA modulus  $N = 336607$  can be distributed. We keep the private key  $d = 280235$  to ourselves, and destroy  $p, q$  and  $\varphi(N)$ , which will not be used anymore.

In order to encrypt a message, it first needs to be translated into a number. For example, one could encode each letter of the alphabet with a number:

$a$	$b$	$c$	$d$	$e$	$f$	$g$	$h$	$i$	$j$	$k$	$l$	$m$
01	02	03	04	05	06	07	08	09	10	11	12	13

<sup>1</sup>In practice, one uses much bigger primes, e.g. primes with 2084 bits, which is more than 600 decimal digits.

<sup>2</sup>One usually chooses a relatively small public key  $e$  which makes the encryption faster, e.g.  $e = 3$  or  $e = 65537 = 2^{16} + 1$ .

### 3 Modular arithmetic

$n$	$o$	$p$	$q$	$r$	$s$	$t$	$u$	$v$	$w$	$x$	$y$	$z$
14	15	16	17	18	19	20	21	22	23	24	25	26

For example, the message “prime” would be encoded as

$$\text{prime} = 16\ 18\ 09\ 13\ 05.$$

Secondly, a message (encoded as a number) may not be larger than  $N$ , and hence must possibly be divided into smaller blocks of size  $< N$  which can then be encrypted individually. Now the *encryption* and *decryption* of a message  $m$  with  $1 \leq m < N$  works as follows:

1. *Encryption*: Compute  $c \equiv m^e \pmod{N}$ , using the public key  $e$ .
2. *Decryption*: Compute  $m \equiv c^d \pmod{N}$ , using the private key  $d$ .

Note that we always take the least residue of  $m^e$  and  $c^d$  modulo  $N$ , that is, we always work with number between 1 and  $N$ . The necessary exponentiation in the encryption and decryption can quickly be computed using the repeated squaring method.

**Proposition 3.5.10.** *The RSA decryption works: if  $c \equiv m^e \pmod{N}$ , then  $m \equiv c^d \pmod{N}$ . Moreover, if  $1 \leq m < N$ , then  $m$  equals the least residue of  $c^d$  modulo  $N$ .*

*Proof.* If  $m$  is coprime to  $N$ , then we can use Euler’s Theorem to compute

$$c^d \equiv (m^e)^d \equiv m^{ed} \equiv m^{1+k\varphi(N)} \equiv m \cdot (m^{\varphi(N)})^k \equiv m \pmod{N},$$

since  $m^{\varphi(N)} \equiv 1 \pmod{N}$  if  $\gcd(m, N) = 1$ . The RSA decryption also works if  $\gcd(m, N) > 1$ , but we leave the verification to the reader.  $\square$

**Remark 3.5.11.** Note that the decryption would not work if we started with a message  $m > N$ : the decryption of the ciphertext  $c \equiv m^e \pmod{N}$  yields the least residue of  $m$  modulo  $N$ , so the original  $m$  and decrypted message  $c^d \pmod{N}$  would differ by some (unknown) multiple of  $N$ .

**Example 3.5.12.** As in the last example, we use the RSA modulus

$$N = 336607$$

and the public and private keys

$$e = 2147, \quad d = 280235.$$

We want to encrypt the message “numbertheory”. It will be encoded as

$$m = 14\ 21\ 13\ 02\ 05\ 18\ 20\ 08\ 05\ 15\ 18\ 25.$$

Since  $m$  is larger than  $N$ , and  $N$  is a 6-digit number, it will be convenient to split  $m$  into blocks of length 6, so we need to encode

$$\begin{aligned} m_1 &= 142113 \\ m_2 &= 020518 \\ m_3 &= 200805 \\ m_4 &= 151825 \end{aligned}$$

### 3.5 Applications to cryptography

Note that we can omit the leading 0 in  $m_2 = 020518$ , since the receiver will see after decrypting  $m_2$  that the number of digits is odd, and hence a leading 0 is missing. Computing  $m_i^d \pmod{N}$ , we obtain the encrypted messages

$$\begin{aligned}c_1 &= 136922 \\c_2 &= 134479 \\c_3 &= 35520 \\c_4 &= 208238\end{aligned}$$

The ciphertexts  $c_1, c_2, c_3, c_4$  can now be sent to the receiver.

In order to decrypt the original message, we need to compute (the least residues of) each  $c_i^d \pmod{N}$ . For example, for  $c_1$  we get

$$c_1^d \equiv 136922^{280235} \equiv 142113 \equiv m_1 \pmod{336607},$$

and putting all decrypted message back together and decoding it, we obtain the original message  $m$  back.

**Remark 3.5.13.** The security of the RSA system depends on two assumptions:

1. An encrypted message  $c$  cannot be decrypted efficiently without the private key  $d$ .
2. The private key  $d$  cannot be computed efficiently from the public key  $e$  and the modulus  $N$ .

The word “efficiently” is crucial here: in principle it is possible to decrypt  $c$  by just trying every possible  $m$  between 1 and  $N$  and checking if  $m^e \equiv c \pmod{N}$ . However, for large  $N$  this becomes unfeasible. Similarly,  $d$  could be computed by factoring  $N = pq$ , computing  $\varphi(N) = (p-1)(q-1)$ , and inverting  $e$  modulo  $\varphi(N)$  using the Euclidean Algorithm. The main assumption in the RSA cryptosystem is that factoring  $N = pq$  becomes computationally impossible for large enough  $N$ .

In order to efficiently compute the private key  $d$ , it would suffice to compute  $\varphi(N)$  from  $N$ , i.e. it seems that factoring  $N$  would not be necessary. However, we will see in the exercise that the factorization of  $N = pq$  can easily be computed if  $\varphi(N)$  is known, so computing  $\varphi(N)$  from  $N$  is computationally equivalent to factoring  $N$ .





## 4 Quadratic forms

### 4.1 Sums of squares

It is clear that not every natural number is the square of a natural number: the squares are  $1, 4, 9, 16, \dots$ , so already 2 is not a square. However, 2 may be written as a *sum of two squares*,

$$2 = 1^2 + 1^2,$$

but it is easy to see that 3 is not the sum of two squares. The classical question which natural numbers can be written as a sum of two squares was first described by Girard in 1625, and further studied by Fermat in 1640. His result, which was proved by Euler in 1752, is as follows.

**Theorem 4.1.1** (Fermat's Two Squares Theorem). *Let  $n$  be a natural number. Then  $n$  can be expressed as a sum of two squares if and only if each prime factor  $p$  of  $n$  with  $p \equiv 3 \pmod{4}$  occurs to an even exponent.*

An important ingredient in the proof is the following simple lemma.

**Lemma 4.1.2** (Diophantus' two squares identity). *For  $a, b, c, d$  in a commutative ring we have*

$$(a^2 + b^2)(c^2 + d^2) = (ac + bd)^2 + (ad - bc)^2.$$

In other words, the product of two natural numbers which are sums of two squares is again a sum of two squares. Hence, it suffices to prove the Two Squares Theorem for  $n = p$  being a prime. Since the Two Squares Theorem is often stated for primes only, we repeat it here:

**Theorem 4.1.3.** *An odd prime  $p$  is a sum of two squares if and only if  $p \equiv 1 \pmod{4}$ .*

This was first proved by Euler, using the method of infinite descent. The proof is completely elementary, but somewhat lengthy. Since we will use this method several times later, we will not do this proof here. We will later see an elegant proof due to Dirichlet, using the theory of binary quadratic forms.

We have seen that not every natural number can be written as a sum of two squares, so the next natural step would be to consider sums of *three* squares. It is easy to find representations for  $1, 2, 3, 4, 5, 6$  as sums of three squares, but for 7 this is not possible. More generally, we have the following characterization.

**Theorem 4.1.4** (Legendre's Three Squares Theorem). *Let  $n$  be a natural number. Then  $n$  can be expressed as a sum of three squares if and only if  $n$  is not of the form  $n = 4^a(8b + 7)$  for non-negative integers  $a, b$ .*

The proof of this result is more difficult than the proof of the Two Squares Theorem, partly due to the fact that there is no 'multiplicativity' result similar to Diophantus' two

#### 4 Quadratic forms

squares identity. We will not discuss the proof in this lecture, but move on to sums of *four* squares. Now  $7 = 1 + 1 + 1 + 4$  is a sum of four squares, and this time we will not find any counter-examples:

**Theorem 4.1.5** (Lagrange's Four Squares Theorem). *Any natural number  $n$  can be written as a sum of four squares.*

As for the Two Squares Theorem, we can reduce to the case of  $n = p$  being a prime by the following lemma.

**Lemma 4.1.6** (Euler's four squares identity). *For  $a_1, \dots, a_4, b_1, \dots, b_4$  in a commutative ring we have*

$$\begin{aligned} (a_1^2 + a_2^2 + a_3^2 + a_4^2)(b_1^2 + b_2^2 + b_3^2 + b_4^2) &= (a_1b_1 - a_2b_2 - a_3b_3 - a_4b_4)^2 \\ &\quad + (a_1b_2 + a_2b_1 + a_3b_4 - a_4b_3)^2 \\ &\quad + (a_1b_3 - a_2b_4 + a_3b_1 + a_4b_2)^2 \\ &\quad + (a_1b_4 + a_2b_3 - a_3b_2 + a_4b_1)^2. \end{aligned}$$

This means that the product of two numbers which are sums of four squares is again a sum of four squares.

*Proof of Theorem 4.1.5.* By Euler's four squares identity, it suffices to prove that every prime  $p$  can be written as a sum of four squares. Since  $2 = 1^2 + 1^2 + 0^2 + 0^2$  we can assume that  $p > 2$  is an odd prime.

We first show that there exists some  $n \in \mathbb{N}$  such that  $np$  can be written as a sum of four squares. It is easy to check that the residues

$$a^2 \pmod{p} \quad \text{for } a \in \{0, \dots, (p-1)/2\}$$

are pairwise distinct. Similarly, the residues

$$-b^2 - 1 \pmod{p} \quad \text{for } b \in \{0, \dots, (p-1)/2\}$$

are pairwise distinct. Since both these sets of residues have  $\frac{p-1}{2} + 1 > p/2$  elements, by the pigeonhole principle there are  $a, b \in \{0, \dots, (p-1)/2\}$  such that

$$a^2 \equiv -b^2 - 1 \pmod{p},$$

or, in other words, there exists some  $n \in \mathbb{N}$  such that

$$a^2 + b^2 + 1^2 + 0^2 = np.$$

Hence  $np$  is a sum of four squares. Note that  $n < p$  by the choice of  $a$  and  $b$ .

If  $n = 1$ , then we are done. Now we assume that  $n > 1$  and construct, from given integers  $x_1, x_2, x_3, x_4$  with

$$x_1^2 + x_2^2 + x_3^2 + x_4^2 = np$$

some new integers  $w_1, w_2, w_3, w_4$  and  $1 \leq r < n$  such that

$$w_1^2 + w_2^2 + w_3^2 + w_4^2 = rp.$$

By repeating this construction until  $r = 1$ , we obtain a representation of  $p$  as a sum of four squares.

To this end, we choose  $y_i$  between  $(-n + 1)/2$  and  $n/2$  (possibly included) with  $y_i \equiv x_i \pmod{n}$ . Then it is easy to show that we have

$$y_1^2 + y_2^2 + y_3^2 + y_4^2 = nr$$

for some natural number  $r$  with  $1 \leq r < n$ . Since  $np$  and  $nr$  are both sums of four squares, Euler's four squares identity shows that  $n^2pr = n^2pr$  is a sum of four squares,

$$n^2pr = z_1^2 + z_2^2 + z_3^2 + z_4^2,$$

where  $z_1, \dots, z_4$  can be written explicitly in terms of the  $x_i$  and  $y_i$ . Plugging in the precise formula for  $z_i$  from Lemma 4.1.6, it is easy to check that  $z_1, \dots, z_4$  are divisibly by  $n$ . Setting  $w_i = z_i/n \in \mathbb{Z}$  we obtain

$$pr = w_1^2 + w_2^2 + w_3^2 + w_4^2,$$

that is,  $pr$  is a sum of four squares. If  $r = 1$ , we are done, and otherwise we continue with the same steps to produce an even smaller  $r$  until we reach  $r = 1$ .  $\square$

We have seen that there are convenient characterizations for numbers which can be written as a sum of two, three, or four squares. But one can even *count* the number of ways in which a number can be written as a sum of squares in these cases. For  $k \in \mathbb{N}$  we let

$$r_k(n) = \#\{(x_1, \dots, x_k) \in \mathbb{Z}^k : x_1^2 + \dots + x_k^2 = n\}$$

be the *sum-of-squares function*. For example, for  $k = 4$  we have:

**Theorem 4.1.7** (Jacobi's Four Squares Theorem). *The number of ways to write a natural number  $n$  as a sum of four squares is given by*

$$r_4(n) = 8\sigma(n) - 32\sigma(n/4) = 8 \sum_{\substack{d|n \\ 4 \nmid d}} d,$$

where  $\sigma(n) = \sum_{d|n} d$  is the *sum-of-divisors function* (and we understand that  $\sigma(n/4) = 0$  unless  $n$  is divisibly by 4), and the sum on the right ranges over all positive divisors  $d$  of  $n$  which are not divisibly by 4.

The standard proof is a beautiful application of the theory of modular forms, but it would lead us too far away here. Since  $d = 1$  always appears in the sum  $\sum_{d|n, 4 \nmid d} d$ , we see that every natural number can be written as a sum of four squares in at least 8 ways, so Jacobi's Theorem implies Lagrange's Theorem. There are similar formulas for  $r_2(n)$ ,  $r_3(n)$ , and  $r_8(n)$ .

## 4.2 Binary quadratic forms

**Definition 4.2.1.** An *integral binary quadratic form* is a homogeneous polynomial of degree two in two variables with integer coefficients,

$$Q(x, y) = ax^2 + bxy + cy^2, \quad (a, b, c \in \mathbb{Z}).$$

The *discriminant* of  $Q$  is defined by

$$D = b^2 - 4ac.$$

We say that  $Q$  is *primitive* if  $\gcd(a, b, c) = 1$ .

## 4 Quadratic forms

For an integral binary quadratic form  $Q(x, y) = ax^2 + bxy + cy^2$  we will often just write  $Q = [a, b, c]$ , and call  $Q$  a *quadratic form* or just a *form*.

**Remark 4.2.2.** The *Gram matrix* of a quadratic form  $Q = [a, b, c]$  is defined by

$$A = \begin{pmatrix} a & b/2 \\ b/2 & c \end{pmatrix}.$$

Then  $Q(x, y)$  can be represented as a matrix product via

$$Q(x, y) = \begin{pmatrix} x & y \end{pmatrix} A \begin{pmatrix} x \\ y \end{pmatrix},$$

and the discriminant of  $Q$  is given by  $-4 \det(A)$ . This representation is useful in computations. By a slight abuse of notation, we will often just write  $Q$  to denote the Gram matrix.

Note that  $b$  and the discriminant  $D = b^2 - 4ac$  have the same parity. Moreover, the discriminant  $D$  of a quadratic form satisfies  $D \equiv 0 \pmod{4}$  or  $D \equiv 1 \pmod{4}$ . Conversely, for any integer  $D \equiv 0, 1 \pmod{4}$  there is at least one quadratic form of discriminant  $D$ :

**Definition 4.2.3.** Let  $D \equiv 0, 1 \pmod{4}$ . The *principal form* of discriminant  $D$  is defined by

$$Q_0 = \begin{cases} [1, 0, -\frac{D}{4}] & \text{if } D \equiv 0 \pmod{4}, \\ [1, 1, \frac{1-D}{4}] & \text{if } D \equiv 1 \pmod{4}. \end{cases}$$

Hence, we will call any integer  $D \equiv 0, 1 \pmod{4}$  a *discriminant*.

**Definition 4.2.4.** We say that  $Q$  *represents* a number  $n \in \mathbb{Z}$  if there exist some  $x, y \in \mathbb{Z}$  with  $Q(x, y) = n$ . Moreover, we say that  $Q$  *properly represents*  $n$  if there exist *coprime*  $x, y \in \mathbb{Z}$  with  $Q(x, y) = n$ .

**Example 4.2.5.** The quadratic form  $Q(x, y) = x^2 + y^2$  is the principal form of discriminant  $-4$ . The question whether  $Q$  represents a number  $n$  is the same as asking whether  $n$  can be written as a sum of two squares.

In general, it is a difficult question to decide which natural numbers are represented by a given quadratic form.

Another basic question is whether a quadratic form can represent both positive and negative integers, or only positive (resp.) negative integers.

**Definition 4.2.6.** A quadratic form  $Q$  is called

- *positive (resp. negative) definite* if  $Q(x, y) > 0$  (resp.  $Q(x, y) < 0$ ) for all  $(x, y) \neq (0, 0)$ .
- *indefinite* if it takes positive and negative values, that is, if there are  $(x, y) \in \mathbb{R}^2$  with  $Q(x, y) > 0$  and  $(x', y') \in \mathbb{R}^2$  with  $Q(x', y') < 0$ .

We have the following handy criterion for definiteness.

**Lemma 4.2.7.** *An integral binary quadratic form  $Q = [a, b, c]$  of discriminant  $D \neq 0$  is*

1. *positive (resp. negative) definite if and only if  $D < 0$  and  $a > 0$  (resp.  $a < 0$ ),*

2. indefinite if and only if  $D > 0$ .

*Proof.* We have

$$4aQ(x, y) = (2ax + by)^2 - Dy^2, \quad (4.2.1)$$

which implies the lemma.  $\square$

**Remark 4.2.8.** The formula in (4.2.1) shows that, if  $D$  is a square, then  $Q(x, y)$  factors over the integers as

$$4aQ(x, y) = (2ax + by - \sqrt{D}y)(2ax + by + \sqrt{D}y),$$

at least if  $a \neq 0$ . However, if  $a = 0$ , then  $Q(x, y) = (bx + c)y$  also factors. In particular, the case of  $D$  being a square (which includes  $D = 0$ ) is not very interesting, so we will assume throughout that  $D$  is not a square. This also implies that  $a, c \neq 0$ .

If  $Q$  is negative definite, then  $-Q$  is positive definite, so we will restrict our attention to positive definite and indefinite forms from now on.

We let

$$\mathrm{SL}_2(\mathbb{Z}) = \left\{ M = \begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix} \in M_2(\mathbb{Z}) : \det(M) = 1 \right\}$$

be the special linear group.

**Definition 4.2.9.** Let  $Q = \begin{pmatrix} a & b/2 \\ b/2 & c \end{pmatrix}$  be (the Gram matrix of) a quadratic form, and let  $M \in \mathrm{SL}_2(\mathbb{Z})$ . We define a new quadratic form by

$$Q \circ M = M^t Q M.$$

Two quadratic forms  $Q, Q'$  are called *equivalent*, written  $Q \sim Q'$ , if there exists a matrix  $M \in \mathrm{SL}_2(\mathbb{Z})$  such that  $Q' = Q \circ M$ .

**Remark 4.2.10.** In group-theoretic terms, the group  $\mathrm{SL}_2(\mathbb{Z})$  *acts* on the set of all quadratic forms by  $Q \circ M$ . Moreover, it is easy to check that the notion of equivalence defines an equivalence relation on the set of all quadratic forms.

If we rather want to work with quadratic forms  $Q(x, y)$  as functions of  $x, y$ , the action of  $M = \begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix} \in \mathrm{SL}_2(\mathbb{Z})$  is given by

$$\left( Q \circ \begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix} \right) (x, y) = Q(\alpha x + \beta y, \gamma x + \delta y).$$

Equivalence of quadratic forms just means that  $Q$  arises from  $Q'$  by an orientation preserving base change  $\mathbb{Z}^2 \rightarrow \mathbb{Z}^2, x \mapsto Mx$ , for some  $M \in \mathrm{SL}_2(\mathbb{Z})$ . This implies the following basic facts:

**Lemma 4.2.11.** *Let  $Q = [a, b, c]$  and  $Q' = [a', b', c']$  be equivalent. Then*

1.  $Q$  and  $Q'$  (properly) represent the same integers.
2.  $Q$  and  $Q'$  have the same discriminant.
3.  $Q$  is positive definite (resp. indefinite) if and only if  $Q'$  is positive definite (resp. indefinite).
4.  $Q$  is primitive if and only if  $Q'$  is primitive.

*In particular, equivalence of quadratic forms defines an equivalence relation on the set of primitive positive definite (resp. indefinite) quadratic forms of a fixed discriminant  $D$ .*

### 4.3 Reduction theory of binary quadratic forms

For a fixed discriminant  $D$  we let

$$\mathcal{Q}_D = \{Q = [a, b, c] : a, b, c \in \mathbb{Z}, b^2 - 4ac = D\}$$

be the set of all integral binary quadratic forms of discriminant  $D$ . Equivalence  $Q \sim Q'$  of quadratic forms gives a partition of  $\mathcal{Q}_D$  into  $\mathrm{SL}_2(\mathbb{Z})$ -equivalence classes. The first main result in the theory, due to Gauss, is the fact that there are only *finitely many* equivalence classes.

**Theorem 4.3.1** (Gauss). *For each fixed discriminant  $D$ , there are only finitely many  $\mathrm{SL}_2(\mathbb{Z})$ -equivalence classes of quadratic forms of discriminant  $D$ .*

*Proof.* We will show that each quadratic form  $[a, b, c]$  of discriminant  $D$  is equivalent to a so-called *weakly reduced form*  $[a', b', c']$ , which means that

$$|b'| \leq |a'| \leq |c'|. \quad (4.3.1)$$

This will imply the stated result, since there are only finitely many weakly reduced forms: we have

$$|D| = |b^2 - 4ac| = |b'^2 - 4a'c'| \geq 3a'^2$$

so there are only finitely many choices for  $a'$ , hence for  $b'$ , and they determine  $c'$  by  $D = b'^2 - 4a'c'$ .

We consider the two  $\mathrm{SL}_2(\mathbb{Z})$ -matrices<sup>1</sup>

$$S = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}, \quad T = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}.$$

Note that  $T^n = \begin{pmatrix} 1 & n \\ 0 & 1 \end{pmatrix}$  for  $n \in \mathbb{Z}$ . Now  $S$  and  $T^n$  act on a quadratic form  $Q = [a, b, c]$  by

$$\begin{aligned} [a, b, c] \circ S &= [c, -b, a], \\ [a, b, c] \circ T^n &= [a, b + 2na, c + nb + n^2a]. \end{aligned}$$

By definition,  $[a, b, c] \circ S$  and  $[a, b, c] \circ T^n$  are equivalent to  $[a, b, c]$ . We will now show that we can apply  $T^n$  and  $S$  repeatedly until we arrive at a weakly reduced form satisfying (4.3.1).

1. First, there is a unique  $n \in \mathbb{Z}$  such that  $-|a| < b + 2na \leq |a|$ , and applying  $T^n$  with this  $n$ , we arrive at a form

$$[a_1, b_1, c_1] = [a, b + 2na, c + nb + n^2a]$$

with  $|b_1| \leq |a_1|$ . Note that the  $a$ -entry does not change in this step.

2. If  $|a_1| \leq |c_1|$ , we are done, and if  $|a_1| > |c_1|$ , we apply  $S$  to get a form

$$[a_2, b_2, c_2] = [c_1, -b_1, a_1]$$

with  $|a_2| < |c_2|$ . Note that the  $a$ -entry now has become smaller in absolute value. However, since the absolute value of the  $b$ -entry did not change in this step, but the absolute value of the  $a$ -entry has become smaller, we might have  $|b| > |a|$  after this step.

---

<sup>1</sup>One can show that  $\mathrm{SL}_2(\mathbb{Z})$  is generated as a group by  $S$  and  $T$ , but we will not need this fact here.

### 4.3 Reduction theory of binary quadratic forms

We can continue with these two steps, and since the  $a$  entry becomes smaller in absolute value every time, after a finite number of iterations we will arrive at a weakly reduced form satisfying (4.3.1).  $\square$

We want to describe the equivalence classes in more detail, and in particular count their precise number. To this end, one has to distinguish between positive definite and indefinite forms. For simplicity, we will only treat the positive definite forms in the lecture.

**Definition 4.3.2.** The number of equivalence classes of primitive positive definite quadratic forms of discriminant  $D < 0$  is called the *class number* of  $D$ , and is denoted by  $h(D)$ .

In order to compute  $h(D)$ , we define a standard set of representatives for the equivalence classes.

**Definition 4.3.3.** A positive definite quadratic form  $[a, b, c]$  of discriminant  $D < 0$  (and  $a, c > 0$ ) is *reduced* if

1.  $|b| \leq a \leq c$ ,
2. If  $|b| = a$  or  $a = c$ , then  $b \geq 0$ .

As in the proof of Theorem 4.3.1, we see that there are only finitely many reduced forms of discriminant  $D$ .

**Theorem 4.3.4.** *Each positive definite quadratic form is equivalent to a unique reduced form. In particular,  $h(D)$  is given by the number of primitive reduced forms of discriminant  $D$ .*

*Proof.* In the proof of Theorem 4.3.1 we have seen that each quadratic form is equivalent to a weakly reduced form  $[a, b, c]$ , that is,  $|b| \leq a \leq c$  (note that  $a, c > 0$  since  $[a, b, c]$  is positive definite). So it remains to show that a weakly reduced form with  $|b| = a$  or  $a = c$ , but  $b < 0$ , is equivalent to another weakly reduced form with  $b \geq 0$ . If  $|b| = a$  but  $b < 0$ , then  $b = -a$ , so

$$[a, b, c] \circ \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} = [a, b + 2a, a + b + c] = [a, -b, c]$$

is reduced. If  $a = c$  but  $b < 0$ , then

$$[a, b, c] \circ \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} = [c, -b, a] = [a, -b, c]$$

is reduced. This shows that any positive definite form is equivalent to a reduced form.

In order to prove uniqueness, suppose that  $[a, b, c]$  and  $[a', b', c']$  are equivalent and reduced. Then there exists some matrix  $\begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix} \in \text{SL}_2(\mathbb{Z})$  such that

$$[a', b', c'] = [a, b, c] \circ \begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix}.$$

A short computation now shows that

$$a' = a\alpha^2 + b\alpha\gamma + c\gamma^2.$$

#### 4 Quadratic forms

Without loss of generality, we can assume that  $a \geq a'$ . Then, using that  $[a, b, c]$  is reduced, we can estimate

$$a \geq a' = a\alpha^2 + b\alpha\gamma + c\gamma^2 \geq a(\alpha^2 + \gamma^2) + b\alpha\gamma \geq a(\alpha^2 + \gamma^2) - a|\alpha\gamma| \geq a|\alpha\gamma|.$$

Since  $\alpha, \gamma$  are integers with  $\gcd(\alpha, \gamma) = 1$ , this leaves only the possibilities  $(\alpha, \gamma) = (\pm 1, 0), (0, \pm 1)$ , and  $(\pm 1, \pm 1)$  (where the signs in the last one are independent).

We only treat the first case, since the others are similar. Suppose that  $\alpha = \pm 1$  and  $\gamma = 0$ . Then  $\begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix} = \pm \begin{pmatrix} 1 & \beta \\ 0 & 1 \end{pmatrix} = \pm T^\beta$ , and we have

$$[a', b', c'] = [a, b, c] \circ \pm T^\beta = [a, b + 2\beta a, c + \beta a + \beta^2 c].$$

This implies  $a = a'$ . If  $|b| \neq a$  then  $|b| \leq a$  and  $|b'| = |b + 2\beta a| \leq a' = a$  together imply  $\beta = 0$ , hence  $[a, b, c] = [a', b', c']$ . If  $|b| = a$ , then  $b \geq 0$  since  $[a, b, c]$  is reduced, so  $|b'| = |b + 2\beta a| \leq a$  is only possible if  $\beta = 0$  or  $\beta = -1$ . But  $\beta = -1$  would mean that  $|b'| = a'$  and  $b' < 0$ , which contradicts  $[a', b', c']$  being reduced. This implies  $\beta = 0$ , so  $[a, b, c] = [a', b', c']$ .  $\square$

**Remark 4.3.5.** The proofs of Theorem 4.3.1 and Theorem 4.3.4 give an algorithmic way to determine the reduced representative of a given quadratic form  $[a, b, c]$ :

1. Replace  $[a, b, c]$  with  $[a, b, c] \circ \begin{pmatrix} 1 & n \\ 0 & 1 \end{pmatrix} = [a, b + 2na, c + nb + n^2 a]$  such that  $-a < b + 2na \leq a$  (take  $n = \lfloor \frac{a-b}{2a} \rfloor$ ).
2. Replace  $[a, b, c]$  with  $[a, b, c] \circ \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} = [c, -b, a]$  if  $a > c$ .
3. Repeat steps 1. and 2. until  $|b| \leq a \leq c$ .
4. If  $|b| = a$  or  $a = c$ , and  $b < 0$ , then replace  $[a, b, c]$  by  $[a, -b, c]$ .

Moreover, the class number  $h(D)$  can be determined by counting the finitely many primitive reduced forms of discriminant  $D$ : for each integer  $a > 0$  with  $a \leq \sqrt{|D|/3}$  and each integer  $b$  with  $|b| \leq a$  such that  $c = \frac{b^2 - D}{4a}$  is a positive integer, check if  $[a, b, c]$  is primitive and reduced.

**Example 4.3.6.** Let us determine the reduced forms of discriminant  $D = -4$ . Then the only possible choice for  $a > 0$  with  $a \leq \sqrt{|D|/3} \leq \sqrt{4/3}$  is  $a = 1$ , and the only possible choices for  $b$  are  $b = 0$  or  $b = \pm 1$ . Now  $b = \pm 1$  does not have the same parity as  $D$  (so  $c = \frac{b^2 - D}{4a}$  is not an integer). This only leaves  $b = 0$  and  $c = \frac{b^2 - D}{4a} = 1$ . Hence, the only reduced form of discriminant  $D = -4$  is the principal form  $[1, 0, 1]$ , and we have class number  $h(-4) = 1$ . In particular, every positive definite quadratic form of discriminant  $-4$  is equivalent to  $[1, 0, 1]$ .

**Remark 4.3.7.** Heegner showed that there are only finitely many discriminants with class number 1, namely

$$D = -3, -4, -7, -8, -11, -12, -16, -19, -27, -28, -43, -67, -163.$$

His proof contained a gap, which was closed by Stark. Hence, this result is known as the Heegner-Stark Theorem. Later, Goldfeld and Gross-Zagier showed that, for each prescribed class number, there are only finitely many discriminants  $D < 0$  with this class number.



**Remark 4.3.8.** Using binary quadratic forms, we can give a short proof of the fact that every prime  $p \equiv 1 \pmod{4}$  can be written as a sum of two squares, due to Dirichlet: Since  $p \equiv 1 \pmod{4}$ , Lemma 3.4.5 tells us that  $-1$  is a square mod  $p$ . Hence, there exists some  $m \in \mathbb{Z}$  such that  $-1 = m^2 + pk$  for some  $k \in \mathbb{Z}$ . Now the binary quadratic form

$$[p, 2m, -k]$$

has discriminant  $-4$  and is positive definite (since  $D = -4 < 0$  and  $p > 0$ ). In the last example, we have seen that  $[1, 0, 1] = x^2 + y^2$  is the unique reduced form of discriminant  $-4$ , hence  $[p, 2m, -k]$  is equivalent to  $[1, 0, 1]$ . Now  $[p, 2m, -k]$  represents  $p$  (just plug in  $(x, y) = (1, 0)$ ), and since equivalent forms represent the same integers,  $[1, 0, 1]$  represents  $p$ , as well. But this means that there are integers  $x, y$  with  $x^2 + y^2 = p$ .

## 4.4 Gauss composition

Diophantus' two squares identity tells us that

$$(x_1^2 + y_1^2)(x_2^2 + y_2^2) = X^2 + Y^2$$

where  $X = x_1x_2 + y_1y_2$  and  $Y = x_1y_2 - y_1x_2$ . In particular, the quadratic form  $x^2 + y^2 = [1, 0, 1]$  represents all the products of numbers represented by  $[1, 0, 1]$ . More generally, we can ask the following question: given two quadratic forms  $Q_1 = [a_1, b_1, c_1]$  and  $Q_2 = [a_2, b_2, c_2]$ , does there exist a third form  $Q_3 = [a_3, b_3, c_3]$  which represents the products of integers represented by  $Q_1$  and  $Q_2$ ? Gauss showed that this is indeed possible. The form  $Q_3$  is called the *Gauss composition* of  $Q_1$  and  $Q_2$ . Moreover, the set of equivalence classes of primitive quadratic forms of a fixed discriminant becomes a finite abelian group under Gauss composition.

Gauss' original definition is quite technical, so it will be easier to work with the *united forms* of Dirichlet.

**Definition 4.4.1.** Two primitive quadratic forms  $[a_1, b_1, c_1]$  and  $[a_2, b_2, c_2]$  of the same discriminant are called *united* if  $\gcd(a_1, a_2, \frac{b_1 + b_2}{2}) = 1$ .

Notice that for forms of the same discriminant,  $\frac{b_1 + b_2}{2}$  is an integer, so  $\gcd(a_1, a_2, \frac{b_1 + b_2}{2})$  is well-defined.

In order to define the Gauss composition of two united forms, we need the following technical lemma.

**Lemma 4.4.2.** *If  $Q_1 = [a_1, b_1, c_1]$  and  $Q_2 = [a_2, b_2, c_2]$  are united forms of discriminant  $D$ , then there exists an integer  $B$  such that*

$$\begin{aligned} B &\equiv b_1 \pmod{2a_1} \\ B &\equiv b_2 \pmod{2a_2} \\ B^2 &\equiv D \pmod{4a_1a_2}. \end{aligned}$$

*In particular,*

$$C = \frac{B^2 - D}{4a_1a_2}$$

*is an integer, and we have the equivalences*

$$\begin{aligned} [a_1, b_1, c_1] &\sim [a_1, B, a_2C], \\ [a_2, b_2, c_2] &\sim [a_2, B, a_1C]. \end{aligned}$$

#### 4 Quadratic forms

*Proof.* Applying  $T^n$  to  $Q_1$  and  $Q_2$  we see that we can change  $b_1$  modulo  $2a_1$  and  $b_2$  modulo  $2a_1$ . Hence, we would like to show that the congruences

$$\begin{aligned} B &\equiv b_1 \pmod{2a_1} \\ B &\equiv b_2 \pmod{2a_2} \end{aligned} \tag{4.4.1}$$

are simultaneously solvable, in such a way that  $C = \frac{B^2 - D}{4a_1a_2}$  is an integer.

The solutions of the first congruence are of the form  $b_1 + 2a_1n$ , and this is a solution for the second congruence if and only if

$$\frac{b_1 - b_2}{2} \equiv -a_1n \pmod{a_2},$$

which has a solution  $n$  if and only if  $d = \gcd(a_1, a_2)$  divides  $(b_1 - b_2)/2$ . From

$$D = b_1^2 - 4a_1c_1 = b_2^2 - 4a_2c_2$$

we see that

$$\left(\frac{b_1 + b_2}{2}\right) \left(\frac{b_1 - b_2}{2}\right) = a_1c_1 - a_2c_2,$$

and since  $d = \gcd(a_1, a_2)$  divides the right-hand side, it must divide the left-hand side. But  $\gcd(a_1, a_2, \frac{b_1 + b_2}{2}) = 1$  since the forms are united, so  $d$  must divide  $(b_1 - b_2)/2$ . This shows that we can find an integer  $B$  satisfying the above two congruences (4.4.1).

Now we show that we can adjust  $B$  such that  $C = \frac{B^2 - D}{4a_1a_2}$  becomes an integer. Let  $k = a_1a_2/d$ . Again by acting with  $T^n$ , we can change  $B$  modulo integer multiples of  $2k$ , and then  $B$  will still satisfy the congruences (4.4.1). Let us write  $B = B_0 + 2kt$  with some fixed choice of  $B_0$  satisfying the congruences, and an integer  $t$  which we would like to choose in such a way that  $B^2 \equiv D \pmod{4a_1a_2}$ . This is equivalent to

$$\frac{D - B_0^2}{4k} \equiv B_0t \pmod{d}.$$

Note that the left-hand side is indeed an integer. Since the two forms are united,  $B_0$  is invertible modulo  $d$ , so

$$t \equiv \left(\frac{D - B_0^2}{4k}\right) B_0^{-1} \pmod{d}$$

determines an appropriate  $B$  for which  $C$  is an integer. □

**Remark 4.4.3.** If  $\gcd(a_1, a_2, \frac{b_1 + b_2}{2}) = 1$ , then there exist integers  $n_1, n_2, n_3$  such that

$$a_1n_1 + a_2n_2 + \frac{b_1 + b_2}{2}n_3 = 1,$$

and we can choose

$$B = a_1b_2n_1 + a_2b_1n_2 + \frac{b_1b_2 + D}{2}n_3.$$

**Definition 4.4.4.** Let  $Q_1 = [a_1, b_1, c_1]$  and  $Q_2 = [a_2, b_2, c_2]$  be united forms of discriminant  $D$ . The *Gauss composition* of  $Q_1$  and  $Q_2$  is defined by

$$Q_1 * Q_2 = [a_1a_2, B, C],$$

where  $B$  and  $C$  are defined as in the last lemma.

One can show that Gauss composition defines a well-defined operation on *equivalence classes* of primitive quadratic forms of the same discriminant. Moreover, this operation is associative and commutative. Unfortunately, the proofs of these facts are rather technical and not very enlightening, so we omit them.

In order to compute the Gauss composition of two primitive forms  $Q_1, Q_2$  we first replace them by equivalent united forms, and then by the equivalent forms  $[a_1, B, a_2C], [a_2, B, a_1, C]$  using Lemma 4.4.2.

**Example 4.4.5.** We consider the two reduced forms

$$Q_1 = [2, 1, 5], \quad Q_2 = [3, 3, 4]$$

of discriminant  $-39$ . Since  $\gcd(2, 3, \frac{1+3}{2}) = 1$ , they are united. We need to determine  $B$  satisfying the above congruences. In our case, this means that

$$\begin{aligned} B &\equiv 1 \pmod{4} \\ B &\equiv 3 \pmod{6} \\ B^2 &\equiv -39 \pmod{24}. \end{aligned}$$

Since  $-39 \equiv 9 \pmod{24}$ , we see that a suitable choice would be  $B = -3$ . Then  $C = \frac{B^2 - D}{4a_1a_2} = \frac{9 + 39}{24} = 2$ , so we have

$$\begin{aligned} [2, 1, 5] &\sim [2, -3, 3 \cdot 2], \\ [3, 3, 4] &\sim [3, -3, 2 \cdot 2]. \end{aligned}$$

Hence, by the definition of Gauss composition, we have

$$[2, 1, 5] * [3, 3, 4] = [2 \cdot 3, -3, 2] = [6, -3, 2].$$

Applying  $S$  and  $T^{-2}$ , we see that

$$[2, 1, 5] * [3, 3, 4] = [6, -3, 2] \sim [2, 3, 6] \sim [2, -1, 5],$$

which is the reduced representative of the Gauss composition  $[2, 1, 5] * [3, 3, 4]$ .

Summarizing, we have the following beautiful result.

**Theorem 4.4.6.** *The set of equivalence classes of primitive (positive definite if  $D < 0$ ) quadratic forms of discriminant  $D$  is a finite abelian group under Gauss composition, called the class group  $\text{Cl}(D)$  of discriminant  $D$ . Moreover, we have:*

1. *The neutral element in  $\text{Cl}(D)$  is given by the class of the principal form,*

$$1_{\text{Cl}(D)} = \begin{cases} \left[ 1, 0, -\frac{D}{4} \right] & \text{if } D \equiv 0 \pmod{4}, \\ \left[ 1, 1, \frac{1-D}{4} \right] & \text{if } D \equiv 1 \pmod{4}. \end{cases}$$

2. *The inverse of  $[a, b, c]$  in  $\text{Cl}(D)$  is given by*

$$[a, b, c]^{-1} = [a, -b, c] \sim [c, b, a].$$

## 4 Quadratic forms

*Proof.* We prove item (1), and leave the second one as an exercise. Moreover, we only do the case  $D \equiv 0 \pmod{4}$ , since the case  $D \equiv 1 \pmod{4}$  is similar. Then  $b$  is even, say  $b = 2n$ , and applying  $T^n = \begin{pmatrix} 1 & n \\ 0 & 1 \end{pmatrix}$  to  $[1, 0, -D/4]$  yields the equivalent form

$$[1, 0, -D/4] \circ T^n = [1, 2n, c'] = [1, b, c'],$$

where  $c' = \frac{b^2 - D}{4}$ . Since  $c = \frac{b^2 - D}{4a}$ , we find  $c' = ac$ , so  $[1, 0, -D/4]$  is equivalent to  $[1, b, ac]$ . Now we can apply the definition of Gauss composition, and obtain

$$[a, b, c] * \left[1, 0, -\frac{D}{4}\right] = [a, b, c] * [1, b, ac] = [a, b, c],$$

as desired. □

## 4.5 Pell's equation and continued fractions

### 4.5.1 Pell's equation

**Definition 4.5.1.** For  $d > 0$  a fixed non-square positive integer, the equation

$$x^2 - dy^2 = 1,$$

is called *Pell's equation*.

Note that the quadratic form  $x^2 - dy^2 = [1, 0, -d]$  has discriminant  $4d > 0$ . Pell's equation always has the *trivial solutions*  $(x, y) = (\pm 1, 0)$ , and if  $d$  is a perfect square, there cannot be any other solutions. We will be interested in its *non-trivial* integral solutions  $(x, y)$  with  $y \neq 0$ . By replacing  $x$  and  $y$  with their negatives, we may focus on *positive* non-trivial solutions  $(x, y) \in \mathbb{N}^2$ . The main goal of this section is to show that all non-trivial solutions to Pell's equation can be generated from a unique *fundamental solution*.

**Definition 4.5.2.** The *fundamental solution*  $(x_1, y_1) \in \mathbb{N}^2$  to Pell's equation  $x^2 - dy^2 = 1$  is the unique non-trivial solution in  $\mathbb{N}^2$  with smallest positive  $x$ -value.

Equivalently, the fundamental solution is the one with smallest positive  $y$ -value, or the solution for which  $x + y\sqrt{d} > 1$  is minimal.

We remark that it is not obvious that  $x^2 - dy^2 = 1$  has any non-trivial solutions, so it is not yet clear that a fundamental solution exists. We will explicitly construct a fundamental solution using continued fractions later. Also note that  $x^2 - dy^2 = 1$  implies that  $\sqrt{d}y < x$ . Hence, we can find the fundamental solution by a “brute-force” search.

We show that *all* positive non-trivial solutions can be generated from the fundamental one.

**Theorem 4.5.3** (Lagrange). *If  $(x_1, y_1)$  is the fundamental solution of Pell's equation  $x^2 - dy^2 = 1$ , where  $d$  is positive and non-square, then every non-trivial solution  $(x, y) \in \mathbb{N}^2$  is given by  $(x_n, y_n)$  where  $x_n + y_n\sqrt{d} = (x_1 + y_1\sqrt{d})^n$  for  $n = 1, 2, 3, \dots$*

**Remark 4.5.4.** The solutions  $(x_n, y_n)$  can be computed recursively via

$$\begin{aligned} x_{n+1} &= x_1x_n + dy_1y_n, \\ y_{n+1} &= x_1y_n + y_1x_n. \end{aligned}$$

#### 4.5 Pell's equation and continued fractions

*Proof of Theorem 4.5.3.* For  $x_n + y_n\sqrt{d} = (x_1 + y_1\sqrt{d})^n$  we compute

$$x_n^2 - dy_n^2 = (x_n + y_n\sqrt{d})(x_n - y_n\sqrt{d}) = (x_1 + y_1\sqrt{d})^n(x_1 - y_1\sqrt{d})^n = (x_1^2 - y_1^2d)^n = 1^n = 1,$$

so  $(x_n, y_n)$  is a solution to Pell's equation.

Conversely, let  $(a, b) \in \mathbb{N}^2$  be a non-trivial solution to  $x^2 - dy^2 = 1$ , and suppose that  $(a, b)$  is not of the form  $(x_n, y_n)$ . We will derive a contradiction by constructing from  $(a, b)$  a new solution  $(u, v) \in \mathbb{N}^2$  to Pell's equation which is smaller than the fundamental solution. Since  $(a, b)$  is not of the form  $(x_n, y_n)$ , there is some positive integer  $k$  such that

$$(x_1 + y_1\sqrt{d})^k < a + b\sqrt{d} < (x_1 + y_1\sqrt{d})^{k+1}.$$

Since  $(x_1 + y_1\sqrt{d})^{-1} = (x_1 - y_1\sqrt{d})$ , dividing by  $(x_1 + y_1\sqrt{d})^k$  we obtain

$$1 < (a + b\sqrt{d})(x_1 - y_1\sqrt{d})^k < x_1 + y_1\sqrt{d}.$$

We put

$$u + v\sqrt{d} = (a + b\sqrt{d})(x_1 - y_1\sqrt{d})^k.$$

Then we have

$$u - v\sqrt{d} = (a - b\sqrt{d})(x_1 + y_1\sqrt{d})^k.$$

We find

$$\begin{aligned} u^2 - dv^2 &= (u + v\sqrt{d})(u - v\sqrt{d}) \\ &= (a + b\sqrt{d})(x_1 - y_1\sqrt{d})^k(a - b\sqrt{d})(x_1 + y_1\sqrt{d})^k \\ &= (a^2 - db^2)(x_1^2 - dy_1^2)^k = 1, \end{aligned}$$

so  $(u, v)$  is again a solution to Pell's equation. By the choice of  $u, v$  we have

$$1 < u + v\sqrt{d} < x_1 + y_1\sqrt{d},$$

so it remains to show that  $u, v > 0$ . Since  $u + v\sqrt{d} > 1$  and  $u - v\sqrt{d} = (u + v\sqrt{d})^{-1}$ , we have

$$0 < u - v\sqrt{d} < 1.$$

Hence

$$2v\sqrt{d} = (u + v\sqrt{d}) - (u - v\sqrt{d}) > 1 - 1 = 0,$$

so  $v > 0$ . Then  $u - v\sqrt{d} > 0$  implies  $u > 0$ . Hence  $(u, v) \in \mathbb{N}^2$  is a solution to Pell's equation which is smaller than the fundamental solution, so we arrive at a contradiction.  $\square$

**Remark 4.5.5.** Non-trivial solutions  $(x, y) \in \mathbb{N}^2$  to Pell's equation  $x^2 - dy^2 = 1$  yield good rational approximations  $\frac{x}{y}$  to  $\sqrt{d}$ . Indeed,  $x^2 - dy^2 = 1$  implies  $x > y\sqrt{d}$  and hence

$$\left| \frac{x}{y} - \sqrt{d} \right| = \frac{1}{y} \cdot \left| \frac{x^2 - dy^2}{x + y\sqrt{d}} \right| = \frac{1}{y} \cdot \left| \frac{1}{x + y\sqrt{d}} \right| < \frac{1}{y(2y\sqrt{d})} < \frac{1}{2y^2},$$

which becomes small for large  $y$ .

### 4.5.2 Finite continued fractions

**Definition 4.5.6.** A *finite continued fraction* of length  $n$  is an iterated sequence of quotients of the form

$$[a_0, a_1, \dots, a_n] = a_0 + \frac{1}{a_1 + \frac{1}{\dots + \frac{1}{a_{n-1} + \frac{1}{a_n}}}},$$

with real numbers  $a_0, a_1, \dots, a_n$ . It is called *simple*, if  $a_0, \dots, a_n \in \mathbb{Z}$  are integers and  $a_1, \dots, a_n \geq 1$ . The numbers  $a_0, \dots, a_n$  are called the *partial quotients* of  $[a_0, a_1, \dots, a_n]$ .

For example, we have

$$\frac{34}{79} = 0 + \frac{1}{\frac{79}{34}} = 0 + \frac{1}{2 + \frac{1}{\frac{34}{11}}} = 0 + \frac{1}{2 + \frac{1}{3 + \frac{1}{11}}} = [0, 2, 3, 11].$$

We have the useful formulas

$$[a_0, \dots, a_n] = a_0 + \frac{1}{[a_1, \dots, a_n]},$$

$$[a_0, \dots, a_n] = [a_0, \dots, a_{n-1} + \frac{1}{a_n}],$$

and it follows from either of the two formulas by induction that every finite simple continued fraction is a rational number. Conversely, we have:

**Theorem 4.5.7.** *Every rational number can be expressed as a finite simple continued fraction. More precisely, the following algorithm computes the finite simple continued fraction expansion of a rational number  $x$ :*

1. Write  $x = a_0 + Z_0$  with  $a_0 \in \mathbb{Z}$  and  $0 \leq Z_0 < 1$ .
2. For  $n \geq 1$ , define  $a_n$  and  $Z_n$  recursively by

$$\frac{1}{Z_{n-1}} = a_n + Z_n, \quad a_n \in \mathbb{N}, \quad 0 \leq Z_n < 1,$$

until  $Z_n = 0$ .

*Proof.* It is clear that the numerator of  $Z_n$  becomes strictly smaller in every step, so it will eventually become 0 and the algorithm stops after finitely many steps.

Moreover, we have

$$x = a_0 + Z_0 = a_0 + \frac{1}{a_1 + Z_1} = a_0 + \frac{1}{a_1 + \frac{1}{a_2 + Z_2}} = \dots = a_0 + \frac{1}{a_1 + \frac{1}{\dots + \frac{1}{a_{n-1} + \frac{1}{a_n}}}},$$

so the algorithm really computes the continued fraction expansion of  $x$ . □

**Remark 4.5.8.** Since

$$[a_0, a_1, \dots, a_{n-1}, a_n] = [a_0, a_1, \dots, a_{n-1}, a_n - 1, 1],$$

the continued fraction expansion of a rational number is not unique. However, if we restrict to simple continued fractions, the above two are the only different representations of the same rational number.

**Definition 4.5.9.** Let  $[a_0, \dots, a_n]$  be a finite continued fraction. The terms

$$\begin{aligned} c_0 &= [a_0] = a_0, \\ c_1 &= [a_0, a_1] = a_0 + \frac{1}{a_1}, \\ c_2 &= [a_0, a_1, a_2] = a_0 + \frac{1}{a_1 + \frac{1}{a_2}}, \\ &\vdots \\ c_k &= [a_0, a_1, \dots, a_k] = a_0 + \frac{1}{a_1 + \frac{1}{\dots + \frac{1}{a_{k-1} + \frac{1}{a_k}}}} \end{aligned}$$

are called the *convergents* of  $[a_0, \dots, a_n]$ .

The convergents can be computed recursively as follows.

**Theorem 4.5.10.** Let  $[a_0, \dots, a_n]$  be a finite continued fraction, and define real numbers  $p_k, q_k$  recursively by

$$\begin{aligned} p_{-1} &= 1, \\ p_0 &= a_0, \\ p_k &= a_k \cdot p_{k-1} + p_{k-2}, \quad (k \geq 1), \end{aligned}$$

and

$$\begin{aligned} q_{-1} &= 0, \\ q_0 &= 1, \\ q_k &= a_k \cdot q_{k-1} + q_{k-2}, \quad (k \geq 1). \end{aligned}$$

Then the  $k$ -th convergent  $c_k = [a_0, \dots, a_k]$  of  $[a_0, \dots, a_n]$  is given by

$$c_k = [a_0, \dots, a_k] = \frac{p_k}{q_k}.$$

The proof is a simple induction over  $k$ , which we leave as an exercise. We emphasize that the theorem holds with real numbers  $a_0, \dots, a_n$ , in which case the  $p_k, q_k$  are also real. Of course, if  $a_0, a_1, \dots, a_n$  are integers, then the  $p_k, q_k$  are integers, as well. The theorem in particular gives the formula

$$[a_0, \dots, a_n, a_{n+1}] = \frac{p_{n+1}}{q_{n+1}} = \frac{p_n a_{n+1} + p_{n-1}}{q_n a_{n+1} + q_{n-1}} \quad (4.5.1)$$

which we will use several times later on.

We collect some fundamental properties of the convergents, all of which are easy to prove and hence are left as an exercise.

## 4 Quadratic forms

**Proposition 4.5.11.** *Let  $x = [a_0, \dots, a_n]$  be a finite continued fraction, and let  $c_k, p_k, q_k$  be as in the last theorem.*

1. *We have the formulas*

$$\begin{aligned} p_k q_{k-1} - p_{k-1} q_k &= (-1)^{k-1} \quad (k \geq 0), \\ \frac{p_k}{q_k} - \frac{p_{k-1}}{q_{k-1}} &= \frac{(-1)^{k-1}}{q_k q_{k-1}} \quad (k \geq 0), \\ p_k q_{k-2} - p_{k-2} q_k &= a_k (-1)^{k-1} \quad (k \geq 1). \end{aligned}$$

2.  $c_0 < c_2 < c_4 < \dots \leq x \leq \dots c_5 < c_3 < c_1$ .

3.  $q_0 < q_1 < \dots < q_n$ .

4. *If  $[a_0, \dots, a_n]$  is simple, then the  $p_k, q_k$  are integers with  $\gcd(p_k, q_k) = 1$ .*

### 4.5.3 Infinite continued fractions

**Definition 4.5.12.** An (infinite) continued fraction is defined by

$$[a_0, a_1, a_2, \dots] = \lim_{n \rightarrow \infty} [a_0, \dots, a_n] = \lim_{n \rightarrow \infty} c_n,$$

where  $a_0$  is an integer, and  $a_1, a_2, a_3, \dots$  is a sequence of positive integers.

We have the suggestive notation

$$[a_0, a_1, a_2, \dots] = a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{a_3 + \dots}}}$$

Also note that, for infinite continued fractions, we always require the convergents  $[a_0, a_1, \dots, a_n]$  to be *simple* finite continued fractions.

**Proposition 4.5.13.** *Under the above conditions on  $a_0, a_1, a_2, \dots$ , the limit*

$$[a_0, a_1, a_2, \dots] = \lim_{n \rightarrow \infty} [a_0, \dots, a_n]$$

*exists, and is greater than any even convergent and smaller than any odd convergent.*

*Proof.* By the fundamental properties of the convergents  $c_n = [a_0, \dots, a_n]$  we have

$$|c_{n+1} - c_n| = \frac{1}{q_n q_{n+1}} < \frac{1}{q_{n+1}^2},$$

which goes to 0 as  $n \rightarrow \infty$  as the  $q_{n+1}$  form a strictly increasing sequence of integers. Hence, the limit exists. Since the even convergents are strictly increasing and the odd convergents are strictly decreasing, the rest of the proposition is clear.  $\square$

**Theorem 4.5.14.** *Every irrational number has an infinite continued fraction expansion, which can be computed using the algorithm from Theorem 4.5.7.*



*Proof.* We have already seen in the proof of Theorem 4.5.7 that

$$x = [a_0, a_1, \dots, a_n, \frac{1}{Z_n}] \tag{4.5.2}$$

for any  $n \geq 0$  (here we allow real partial quotients for the moment). We consider the convergents  $\frac{p_n}{q_n} = [a_0, a_1, \dots, a_n]$ , and claim that

$$\left| x - \frac{p_n}{q_n} \right| = \frac{1}{q_n(q_n/Z_n + q_{n-1})}. \tag{4.5.3}$$

Indeed, by (4.5.2) and (4.5.1) we have

$$x = \frac{p_n/Z_n + p_{n-1}}{q_n/Z_n + q_{n-1}},$$

so we can compute

$$\begin{aligned} x - \frac{p_n}{q_n} &= \frac{p_n/Z_n + p_{n-1}}{q_n/Z_n + q_{n-1}} - \frac{p_n}{q_n} \\ &= \frac{p_{n-1}q_n - p_nq_{n-1}}{q_n(q_n/Z_n + q_{n-1})} \\ &= \frac{(-1)^n}{q_n(q_n/Z_n + q_{n-1})}. \end{aligned}$$

Now for  $n \geq 1$ ,  $a_n$  and  $q_n$  are positive integers, and we have  $1/Z_n = a_{n+1} + Z_{n+1} \geq a_{n+1}$ , so we can estimate

$$\left| x - \frac{p_n}{q_n} \right| < \frac{1}{q_n(q_n a_{n+1} + q_{n-1})} = \frac{1}{q_n q_{n+1}}. \tag{4.5.4}$$

Since the  $q_n$  form a strictly increasing sequence of integers, the right-hand side goes to 0 as  $n \rightarrow \infty$ .  $\square$

**Example 4.5.15.** Let us compute the continued fraction expansion of  $\sqrt{3} \sim 1.732$ . We first write

$$\sqrt{3} = \underbrace{1}_{a_0} + \underbrace{(\sqrt{3} - 1)}_{Z_0}.$$

Now

$$\frac{1}{Z_0} = \frac{1}{\sqrt{3} - 1} = \frac{\sqrt{3} + 1}{(\sqrt{3} - 1)(\sqrt{3} + 1)} = \frac{\sqrt{3} + 1}{2} = \underbrace{1}_{a_1} + \underbrace{\frac{\sqrt{3} - 1}{2}}_{Z_1}.$$

And once more:

$$\frac{1}{Z_1} = \frac{2}{\sqrt{3} - 1} = \sqrt{3} + 1 = \underbrace{2}_{a_2} + \underbrace{\sqrt{3} - 1}_{Z_2}.$$

We see that  $Z_2 = Z_0$ , so the algorithm will repeat itself from now on, e.g. we will get  $Z_3 = Z_1, Z_4 = Z_0, \dots$  and  $a_3 = a_1, a_4 = a_2, a_5 = a_1, \dots$  and so on. Hence, we obtain the *periodic* continued fraction expansion

$$\sqrt{3} = [1, 1, 2, 1, 2, 1, 2, 1, 2, \dots] = [1, \overline{1, 2}].$$

In the next section we will determine precisely the convergents of  $\sqrt{d}$  which give solutions to Pell's equation.

#### 4.5.4 Solutions to Pell's equation via continued fractions

**Definition 4.5.16.** A continued fraction  $[a_0, a_1, a_2, \dots]$  is called *periodic* if the sequence of partial quotients repeats from some point on, that is, if there exist some index  $k_0$  and some fixed natural number  $\ell$  such that  $a_{k+\ell} = a_k$  for all  $k \geq k_0$ . In this case,  $\ell$  is called the *length* of the period, and we write the continued fraction as

$$[a_0, a_1, \dots, a_{k_0-1}, \overline{a_{k_0}, \dots, a_{k_0+\ell-1}}].$$

A periodic continued fraction is called *purely periodic* if  $k_0 = 0$ , that is, if it is of the form

$$[\overline{a_0, a_1, \dots, a_{\ell-1}}]$$

for some natural number  $\ell$ .

We have seen that  $\sqrt{3} = [1, \overline{1, 2}]$  has a periodic continued fraction expansion of length 2.

**Definition 4.5.17.** We call an irrational real number  $w$  a *quadratic irrational* if it is the root of a quadratic polynomial  $ax^2 + bx + c$  with *integer* coefficients  $a, b, c$ . The other root of this polynomial is called the *conjugate* of  $w$ , and is denoted by  $w'$ .

**Remark 4.5.18.** If  $D = b^2 - 4ac > 0$  is not a square, then the pair of conjugate quadratic irrationalities satisfying  $ax^2 + bx + c = 0$  is given by

$$w = \frac{-b + \sqrt{D}}{2a}, \quad w' = \frac{-b - \sqrt{D}}{2a}.$$

In particular, an irrational number  $w$  is a quadratic irrational if and only if it is of the shape  $w = p + q\sqrt{D}$  for some rational numbers  $p, q \in \mathbb{Q}$  with  $q \neq 0$ , and some non-square positive integer  $D$ .

**Theorem 4.5.19.** *An irrational number  $w$  is a quadratic irrational if and only if it has a periodic continued fraction expansion.*

*Proof.* First suppose that  $w$  has a periodic continued fraction expansion,

$$w = [a_0, \dots, a_{k-1}, \overline{a_k, \dots, a_{k+\ell-1}}].$$

We want to show that  $w$  is a quadratic irrational. To this end, we first show that the irrational number

$$\alpha = [\overline{a_k, \dots, a_{k+\ell-1}}]$$

with purely periodic continued fraction expansion is a quadratic irrational. We can write  $\alpha$  as a finite (non-simple) continued fraction

$$\alpha = [a_k, \dots, a_{k+\ell-1}, \alpha].$$

By (4.5.1), we have

$$\alpha = \frac{p\alpha + p'}{q\alpha + q'},$$

where  $p, p', q, q'$  are the rational numbers defined by the convergents  $\frac{p}{q} = [a_k, \dots, a_{k+\ell-1}]$  and  $\frac{p'}{q'} = [a_k, \dots, a_{k+\ell-2}]$ . From this, we get a rational quadratic equation

$$q\alpha^2 + (q' - p)\alpha - p' = 0,$$

#### 4.5 Pell's equation and continued fractions

and clearing the denominators of  $q, (q' - p), p'$ , we obtain an integral quadratic equation for  $\alpha$ , which shows that  $\alpha$  is a quadratic irrational. Using

$$w = [a_0, \dots, a_{k-1}, \alpha],$$

together with the simple fact that sums and quotients of quadratic irrationals of the same discriminant  $D$  are again quadratic irrationals of discriminant  $D$ , we find that  $w$  is a quadratic irrational.

For the converse, we only give the rough idea, since the proof is rather long and technical. Let  $w$  be a quadratic irrational, and a root of  $ax^2 + bx + c = 0$ . Put  $D = b^2 - 4ac$ , which is not a square as  $w$  is irrational. For simplicity, let us assume that  $w$  is the root

$$w = \frac{-b + \sqrt{D}}{2a}.$$

Now we apply the algorithm from Theorem 4.5.7 to compute the continued fraction expansion of  $w$ . Then  $Z_0$  is given by

$$Z_0 = \frac{-b + \sqrt{D}}{2a} + n = \frac{-b + 2an + \sqrt{D}}{2a}$$

for some  $n \in \mathbb{Z}$  which is chosen such that  $0 < Z_0 < 1$ . In particular,  $Z_0$  is of the form

$$Z_0 = \frac{P + \sqrt{D}}{Q} \tag{4.5.5}$$

with integers  $P, Q$  such that  $Q \mid (P^2 - D)$ . Then

$$\frac{1}{Z_0} = \frac{Q}{P + \sqrt{D}} = \frac{Q(P - \sqrt{D})}{P^2 - D}.$$

Since  $Q \mid (P^2 - D)$ , we see that  $Z_1$  will also be of the form (4.5.5). By the same argument, every  $Z_n$  will be of the form (4.5.5), that is, we can write

$$Z_n = \frac{P_n + \sqrt{D}}{Q_n}$$

with integers  $P_n, Q_n$  such that  $Q_n \mid (P_n^2 - D)$ . Now one can use the fundamental properties of the convergents to show that there are actually only finitely many choices for  $P_n, Q_n$ . This is the technical part, which we omit here. Then, by the Pigeonhole principle there exists some  $m$  such that  $Z_n = Z_m$ , so the continued fraction expansion of  $w$  is periodic.  $\square$

**Definition 4.5.20.** A quadratic irrational  $w$  is called *reduced* if  $w > 0$  and  $-1 < w' < 0$ .

We will state the following beautiful result without proof.

**Theorem 4.5.21** (Galois). *A quadratic irrational  $w$  is reduced if and only if it has a purely periodic continued fraction expansion.*

We want to use continued fractions to solve Pell's equation  $x^2 - dy^2 = 1$ , and thereby find good rational approximations to  $\sqrt{d}$ . To this end, we need the shape of the continued fractional expansion of  $\sqrt{d}$ .

#### 4 Quadratic forms

**Proposition 4.5.22.** *Let  $d > 0$  be a non-square natural number. Then  $\sqrt{d}$  has a continued fraction expansion of the form*

$$\sqrt{d} = [a_0, \overline{a_1, a_2, \dots, a_{n-1}, 2a_0}]$$

for some positive integers  $a_0, a_1, a_2, \dots, a_{n-1}$ .

*Proof.* Put

$$\begin{aligned}\alpha &= [\sqrt{d}] - \sqrt{d}, \\ \beta &= [\sqrt{d}] + \sqrt{d}.\end{aligned}$$

Then we have

$$-1 < \alpha < 0, \quad \beta > 1,$$

and

$$(x - \alpha)(x - \beta) = x^2 - (\alpha + \beta)x + \alpha\beta = x^2 + 2[\sqrt{d}]x + [\sqrt{d}]^2 - d,$$

so  $\alpha, \beta$  are conjugate quadratic irrationalities, and  $\beta$  is reduced. Hence, by Galois' Theorem,  $\beta$  has a purely periodic continued fraction expansion

$$\beta = [\overline{a_0, \dots, a_{n-1}}] = [a_0, \overline{a_1, \dots, a_{n-1}, a_0}],$$

where  $a_0 = [\beta] = 2[\sqrt{d}]$ . This implies

$$\sqrt{d} = \beta - [\sqrt{d}] = [2[\sqrt{d}], \overline{a_1, \dots, a_{n-1}, 2[\sqrt{d}]}] - [\sqrt{d}] = [[\sqrt{d}], \overline{a_1, \dots, a_{n-1}, 2[\sqrt{d}]}].$$

This gives the stated continued fraction expansion. □

**Theorem 4.5.23.** *Let  $d > 0$  be a non-square natural number, and let*

$$\sqrt{d} = [a_0, \overline{a_1, a_2, \dots, a_{n-1}, 2a_0}]$$

be the continued fraction expansion of  $\sqrt{d}$ , such that the period  $n$  is minimal.

1. If  $n$  is even, put  $\frac{x}{y} = \frac{p_{n-1}}{q_{n-1}} = [a_0, a_1, \dots, a_{n-1}]$ ,
2. If  $n$  is odd, put  $\frac{x}{y} = \frac{p_{2n-1}}{q_{2n-1}} = [a_0, a_1, \dots, a_{2n-1}]$ .

Then  $(x, y)$  is the fundamental solution to Pell's equation  $x^2 - dy^2 = 1$ .

*Proof.* We can write  $\sqrt{d}$  as a finite continued fraction (involving real numbers) as

$$\sqrt{d} = [a_0, a_1, \dots, a_{n-1}, a_0 + \sqrt{d}].$$

By Theorem 4.5.10, we have

$$\sqrt{d} = \frac{(a_0 + \sqrt{d})p_{n-1} + p_{n-2}}{(a_0 + \sqrt{d})q_{n-1} + q_{n-2}}.$$

Clearing the denominator, we get

$$(a_0q_{n-1} + q_{n-2} - p_{n-1})\sqrt{d} = a_0p_{n-1} + p_{n-2} - q_{n-1}d,$$

which is only possible if both sides are 0, that is,

$$\begin{aligned} p_{n-1} &= a_0 q_{n-1} + q_{n-2} \\ q_{n-1} d &= a_0 p_{n-1} + p_{n-2}. \end{aligned}$$

Hence, we find

$$\begin{aligned} p_{n-1}^2 - d q_{n-1}^2 &= p_{n-1}(a_0 q_{n-1} + q_{n-2}) - q_{n-1}(a_0 p_{n-1} + p_{n-2}) \\ &= p_{n-1} q_{n-2} - q_{n-1} p_{n-2} = (-1)^{n-2}, \end{aligned}$$

where we again used Theorem 4.5.10. For even  $n$ , this shows that  $(p_{n-1}, q_{n-1})$  solve Pell's equation.

For odd  $n$ , we can do the same computation if we write the continued fraction expansion of  $\sqrt{d}$  as

$$\sqrt{d} = [a_0, a_1, \dots, a_{2n-1}, a_0 + \sqrt{d}],$$

(that is, take the periodic part twice) in order to see that then  $(p_{2n-1}, q_{2n-1})$  solves Pell's equation.

We will not prove here the fact that the above solutions are indeed fundamental.  $\square$

**Example 4.5.24.** We have  $\sqrt{3} = [1, \overline{1, 2}]$ , so the period is  $n = 2$ . Since  $n$  is even, we consider the convergent  $c_{n-1} = c_1$ , given by

$$c_1 = [1, 1] = 1 + \frac{1}{1} = 2 = \frac{2}{1}.$$

Indeed,  $(2, 1)$  is the fundamental solution to Pell's equation  $x^2 - 3y^2 = 1$ .

## 4.6 Congruent numbers

### 4.6.1 Pythagorean triples

The Pythagorean Theorem tells us that for any right triangle with legs  $a, b$  and hypotenuse  $c$  we have

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

As usual, we are interested in describing the rational or integral solutions to such Diophantine equations, which leads to the following definition.

**Definition 4.6.1.** A *Pythagorean triple*  $(a, b, c)$  consists of natural numbers with  $a^2 + b^2 = c^2$ .

In other words, the Pythagorean triples are precisely the triples of natural numbers which appear as side lengths of right triangles. The simplest example is the Pythagorean triple  $(3, 4, 5)$ .

A Pythagorean triple  $(a, b, c)$  is called *primitive* if  $a, b, c$  are coprime (which implies that they are also pairwise coprime). Note that, for a primitive triple,  $a$  or  $b$  must have opposite parity, and by switching  $a$  and  $b$  we can always assume that  $a$  is odd.

The primitive Pythagorean triples can be effectively parametrized as follows.

#### 4 Quadratic forms

**Theorem 4.6.2.** For two natural numbers  $m > n$  let

$$a = m^2 - n^2, \quad b = 2mn, \quad c = m^2 + n^2.$$

Then  $(a, b, c)$  is a Pythagorean triple, which is primitive if and only if  $m$  and  $n$  are coprime and of different parity.

Conversely, every primitive Pythagorean triple with odd  $a$  arises in this way from a unique pair of coprime natural numbers  $m > n$  of different parity.

**Corollary 4.6.3.** There are infinitely many primitive Pythagorean triples.

*Proof of Theorem 4.6.2.* We have

$$a^2 + b^2 = (m^2 - n^2)^2 + (2mn)^2 = (m^2 + n^2)^2 = c^2,$$

so the triple  $(a, b, c)$  is Pythagorean.

Suppose that  $(a, b, c)$  arises from  $m$  and  $n$  as above. Let us assume that  $(a, b, c)$  is not primitive. Let  $p$  be a prime dividing  $\gcd(a, b, c)$ . If  $p = 2$ , then  $a = m^2 - n^2$  implies that  $m$  and  $n$  have the same parity. If  $p > 2$ , then  $p$  dividing  $b = 2mn$  implies that  $p$  divides either  $m$  or  $n$ . Since  $p$  also divides  $a = m^2 - n^2$ , it must divide  $m$  and  $n$ , so  $m$  and  $n$  are not coprime. Conversely, if  $m$  and  $n$  are not coprime, then  $\gcd(m, n)$  divides  $\gcd(a, b, c)$ , and if  $m, n$  have the same parity, then  $a, b, c$  are all even, so in both cases  $(a, b, c)$  is not primitive.

We show that any primitive Pythagorean triple  $(a, b, c)$  with odd  $a$  arises from coprime  $m > n$  of different parity. We write

$$\frac{m}{n} = \frac{c+a}{b}$$

in lowest terms, that is,  $\gcd(m, n) = 1$ . Notice that  $\frac{c+a}{b} > 1$  since the sum of any two sides of a triangle must be bigger than the third side. This implies  $m > n$ .

Writing  $b^2 = c^2 - a^2 = (c-a)(c+a)$  we see that

$$\frac{n}{m} = \frac{c-a}{b}.$$

This implies

$$\frac{c}{b} = \frac{1}{2} \left( \frac{m}{n} + \frac{n}{m} \right) = \frac{m^2 + n^2}{2mn} \quad \text{and} \quad \frac{a}{b} = \frac{1}{2} \left( \frac{m}{n} - \frac{n}{m} \right) = \frac{m^2 - n^2}{2mn}.$$

Since  $(a, b, c)$  is primitive,  $a, b, c$  are pairwise coprime. Since  $m$  and  $n$  are coprime, they cannot both be even. If they would both be odd, then the numerator of  $\frac{a}{b} = \frac{m^2 - n^2}{2mn}$  would be divisibly by 4, but the denominator would only be divisibly by 2. But this would imply that  $a$  is even, but we assumed that  $a$  is odd. Hence one of  $m, n$  is even and the other one is odd. This also means that the numerator of  $\frac{m^2 \pm n^2}{2mn}$  is odd. Moreover, it is easy to see that  $\frac{m^2 \pm n^2}{2mn}$  is already reduced: any odd prime dividing the denominator divides  $m$  or  $n$ , but not both, so it cannot divide  $m^2 \pm n^2$ . Since  $\frac{c}{b}$  and  $\frac{a}{b}$  are also reduced, we find

$$a = m^2 - n^2, \quad b = 2mn, \quad c = m^2 + n^2$$

with  $m, n$  coprime and of opposite parity. By adding and subtracting  $a$  and  $c$  it is easy to see that  $m$  and  $n$  are uniquely determined by  $a$  and  $c$ .  $\square$

We have seen that the equation  $a^2 + b^2 = c^2$  has infinitely many solutions in natural numbers which can be described quite explicitly. A natural follow-up question would be to describe the integer solutions of the equation

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

for  $n \geq 3$ . There are some trivial solutions if  $a, b$ , or  $c$  equals 0. However, Fermat proved in 1673 that for  $n = 4$  there are no non-trivial integer solutions, and he conjectured that there will never be a non-trivial solution for  $n \geq 3$ . He claimed to have a “marvelous” proof, which may be doubted in the light of the advanced tools used in the modern proof of Fermat’s conjecture:

**Theorem 4.6.4** (Fermat’s Last Theorem). *For  $n \geq 3$  the equation  $a^n + b^n = c^n$  has no integer solutions  $a, b, c$  with  $abc \neq 0$ .*

Using the method of infinite descent, Euler proved the case  $n = 3$  in 1753. Important special cases of Fermat’s Last Theorem had been proved by Ernst Kummer and Sophie Germain, but the complete proof for every  $n \geq 3$  was only accomplished in 1995 by Andrew Wiles and Richard Taylor, building on fundamental work of many others, e.g. Frey, Ribet and Serre. The proof uses many deep results from the theory of elliptic curves and modular forms.

We will discuss the elementary proof of the case  $n = 4$  in the next section.

#### 4.6.2 Congruent numbers

**Definition 4.6.5.** A natural number  $n$  is called a *congruent number* if it is the area of a right triangle with rational side lengths.

The simplest example is the right triangle with side lengths 3, 4, 5 and area 6, so 6 is a congruent number. Fibonacci found suitable triangles for 5 and 7, showing that they are congruent numbers.

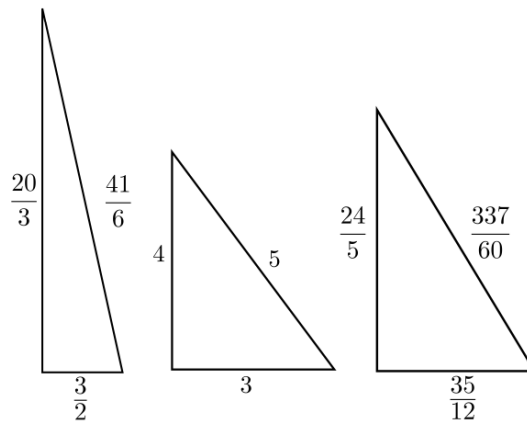


Figure 4.1: Rational right triangles with respective areas 5, 6, 7

However, already for relatively small congruent numbers  $n$ , the side lengths of suitable rational right triangles are typically rational numbers with large numerators and denominator.

## 4 Quadratic forms

A famous example due to Zagier is the ‘simplest’ rational right triangle for the congruent number  $n = 157$ .

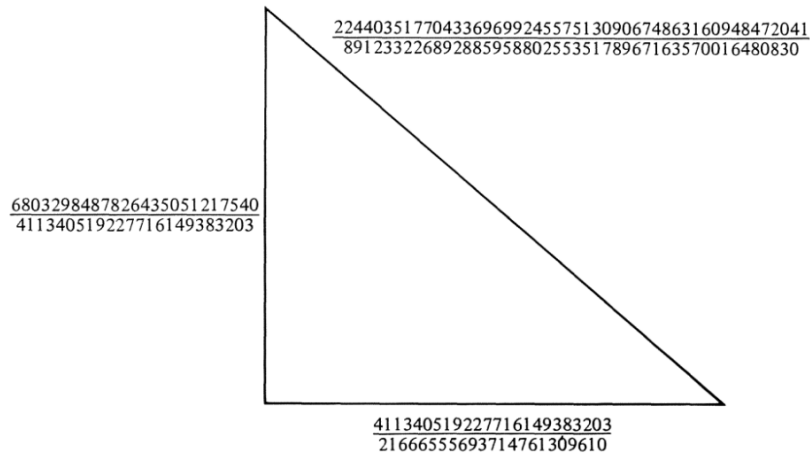


Figure 4.2: Zagier’s rational right triangle with area 157

In particular, finding a suitable right triangle for a given number  $n$  can be quite challenging. However, showing that such a triangle may not exist, seems to be even more difficult.

**Remark 4.6.6.** By rescaling the sides  $a, b, c$  by a common factor  $d \in \mathbb{Q}$ , we see that  $n$  is congruent if and only if  $nd^2$  is congruent, for any  $d \in \mathbb{Q}$  for which  $nd^2$  is a natural number. For example, 1 is congruent if and only if each square of a natural number is congruent. In particular, if we want to investigate congruent numbers, it suffices to consider only square-free natural numbers.

In more algebraic terms,  $n$  is congruent if the equations

$$\begin{aligned} a^2 + b^2 &= c^2 \\ \frac{ab}{2} &= n \end{aligned}$$

have a rational solution  $a, b, c \in \mathbb{Q}$ . In particular, congruent numbers are closely related to Pythagorean triples: any Pythagorean triple  $(a, b, c) \in \mathbb{Z}^3$  yields the congruent number  $n = \frac{ab}{2}$ . Hence, by enumerating all primitive Pythagorean triples  $(a, b, c)$  using their parametrization in terms of  $m, n$  as in Lemma 4.6.2, and taking the square-free part of  $n = \frac{ab}{2}$ , we obtain a simple (yet very inefficient) algorithm that lists square-free congruent numbers together with suitable right triangles. However, if a number  $n$  did not yet appear in this list, we do not know if the reason is that  $n$  is not congruent or that we did not wait long enough.

**Theorem 4.6.7** (Fermat). *1 is not a congruent number.*

**Corollary 4.6.8.** *No square of a natural number is a congruent number.*

*Proof of Theorem 4.6.7.* The proof uses Fermat’s method of infinite descent. The basic idea is to take a ‘minimal’ rational triangle for 1, and from this construct an even ‘smaller’ triangle, which gives a contradiction.



Let's suppose that 1 is congruent, so there exists a rational right triangle with sides  $a, b, c \in \mathbb{Q}$  and area 1. Multiplying  $a, b, c$  by their common denominator, we obtain a Pythagorean triple, and the area of the corresponding right triangle is an integer square. We let  $k^2$  be the smallest area that arises in this way (that is, as the area of a right triangle with *integer* sides), and we let  $(a, b, c)$  be the corresponding Pythagorean triple with  $k^2 = \frac{ab}{2}$ . Since  $k^2$  is minimal,  $(a, b, c)$  must be primitive. Hence (after switching  $a$  and  $b$  if  $a$  is even), by Lemma 4.6.2 there are coprime  $m > n$  of different parity such that

$$a = m^2 - n^2, \quad b = 2mn, \quad c = m^2 + n^2.$$

In terms of  $m$  and  $n$  the area of our right triangle is given by

$$k^2 = \frac{ab}{2} = (m^2 - n^2)mn = (m - n)(m + n)mn.$$

Since  $m, n, m - n$ , and  $m + n$  are pairwise coprime, they must all be integer squares,

$$m = x^2, \quad n = y^2, \quad m + n = u^2, \quad m - n = v^2,$$

for some  $x, y, u, v \in \mathbb{N}$ . We obtain

$$(u + v)^2 + (u - v)^2 = 2u^2 + 2v^2 = 4m = (2x)^2.$$

Hence the triple  $(u + v, u - v, 2x)$  is a Pythagorean triple, whose corresponding area is a square,

$$\frac{(u + v)(u - v)}{2} = \frac{u^2 - v^2}{2} = n = y^2,$$

which is smaller than  $k^2$  since  $n$  is a proper divisor of  $k^2$ . This is a contradiction to the minimality of  $k^2$ , hence 1 cannot be a congruent number.  $\square$

**Remark 4.6.9.** Using the same idea, one can show that 2 and 3 (and hence  $2n^2$  and  $3n^2$  for  $n \in \mathbb{N}$ ) are not congruent.

As a corollary to the fact that 2 is not a congruent number, we obtain a special case of Fermat's last theorem:

**Corollary 4.6.10.** *The equation  $x^4 + y^4 = z^4$  has no integer solutions with  $xyz \neq 0$ .*

*Proof.* Suppose that there are integers  $x, y, z$  with  $x^4 + y^4 = z^4$  and  $xyz \neq 0$ . Then the triangle with side lengths  $(a, b, c) = (x^2, y^2, z^2)$  is a rational right triangle with area  $x^2y^2/2$ . Note that at least one of  $x, y$  must be even (consider the equation  $x^4 + y^4 = z^4$  modulo 4 to see this), say  $x = 2r$ . Now our rational right triangle has area  $x^2y^2/2 = 2(r^2y^2)$ , which is impossible no square multiple of 2 is a congruent number.  $\square$

Next, we would like to relate congruent numbers to another classical number theoretical problem. Consider the 3 integer squares 1, 25, 49. They have common difference  $25 - 1 = 49 - 25 = 24$ . More generally, given a natural number  $n$ , we may ask whether there exists a rational number  $s$  such that  $s^2 - n$  and  $s^2 + n$  are both rational squares. In this case, we say that the three squares  $r^2 = s^2 - n, s^2, t^2 = s^2 + n$  form a *3-term arithmetic progression of rational squares with common difference  $n$* .

#### 4 Quadratic forms

**Proposition 4.6.11.** *A natural number  $n$  is congruent if and only if there exists a 3-term arithmetic progression of rational squares with common difference  $n$ . More precisely, there is a one-to-one correspondence between the sets*

$$\{(a, b, c) \in \mathbb{Q}^3 : a^2 + b^2 = c^2, \frac{ab}{2} = n\} \longleftrightarrow \{(r, s, t) \in \mathbb{Q}^3 : s^2 - r^2 = n, t^2 - s^2 = n\}$$

given by

$$(a, b, c) \mapsto \left( \frac{b-a}{2}, \frac{c}{2}, \frac{b+a}{2} \right), \quad (r, s, t) \mapsto (t-r, t+r, 2s).$$

We leave the proof as an exercise.

**Proposition 4.6.12.** *A natural number  $n$  is congruent if and only if the equation*

$$y^2 = x^3 - n^2x$$

has a rational solution  $(x, y) \in \mathbb{Q}^2$  with  $y \neq 0$ . More precisely, there is a one-to-one correspondence between the sets

$$\{(a, b, c) \in \mathbb{Q}^3 : a^2 + b^2 = c^2, \frac{ab}{2} = n\} \longleftrightarrow \{(x, y) \in \mathbb{Q}^2 : y^2 = x^3 - n^2x, y \neq 0\}$$

given by

$$(a, b, c) \mapsto \left( \frac{n(a+c)}{b}, \frac{2n^2(a+c)}{b^2} \right), \quad (x, y) \mapsto \left( \frac{x^2 - n^2}{y}, \frac{2nx}{y}, \frac{x^2 + n^2}{y} \right).$$

Again, the proof is a nice exercise. The equation  $E_n : y^2 = x^3 - n^2x$  defines a smooth curve in  $\mathbb{R}^2$ , and is an example of an *elliptic curve*. These curves typically look as follows.

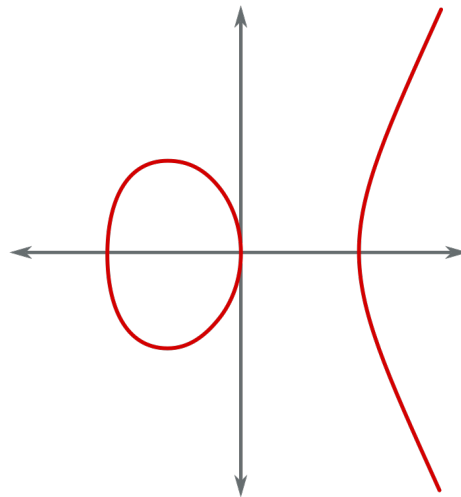


Figure 4.3: The elliptic curve  $E_1 : y^2 = x^3 - x$

The fact that 1 is not congruent is equivalent to the fact that the above curve has no rational point apart from the three points with  $y = 0$ .

Elliptic curves lie at the heart of modern number theory and are the main players in some of the most important breakthroughs and conjectures of the last decades, e.g. in the solution of Fermat's last theorem due to Andrew Wiles, and in the Birch and Swinnerton-Dyer conjecture. Using some deep results about elliptic curves and modular forms, Tunnell was able to find the following handy criterion for congruent numbers.

**Theorem 4.6.13** (Tunnell 1983). *Let  $n$  be a square-free natural number. Consider the numbers*

$$\begin{aligned} A(n) &= \#\{(x, y, z) \in \mathbb{Z}^3 : 2x^2 + y^2 + 8z^2 = n\}, \\ B(n) &= \#\{(x, y, z) \in \mathbb{Z}^3 : 2x^2 + y^2 + 32z^2 = n\}, \\ C(n) &= \#\{(x, y, z) \in \mathbb{Z}^3 : 8x^2 + 2y^2 + 16z^2 = n\}, \\ D(n) &= \#\{(x, y, z) \in \mathbb{Z}^3 : 8x^2 + 2y^2 + 64z^2 = n\}. \end{aligned}$$

Then the following are true.

1. If  $n$  is an odd congruent number, then  $A(n) = 2B(n)$ .
2. If  $n$  is an even congruent number, then  $C(n) = 2D(n)$ .

Moreover, if the weak Birch and Swinnerton-Dyer conjecture holds for the elliptic curve  $E_n : y^2 = x^3 - n^2x$ , then the above coefficients are also sufficient to conclude that  $n$  is congruent.

**Remark 4.6.14.** The four sets in Tunnell's Theorem are certainly finite, since we have the trivial bound  $x, y, z \leq \sqrt{n}$ . We may use the theorem to quickly show that a given number is not congruent. For example, for  $n = 1$  we have  $A(1) = 2$  and  $B(1) = 2$  (the only solutions being  $(0, \pm 1, 0)$ ), so 1 is not congruent. On the other hand, for the congruent number  $n = 5$  we have  $A(5) = 0 = 2B(5)$ , but we cannot conclude from the theorem that 5 is congruent. However, if the condition of Tunnell's Theorem is fulfilled for some  $n$ , it gives a strong hint that  $n$  should be congruent, so we can start searching for a suitable triangle.

There are several open problems about congruent numbers, for example:

**Conjecture 4.6.15.** *Every natural number  $n \equiv 5, 6, 7 \pmod{8}$  is a congruent number.*

This would follow from Tunnell's Theorem if the weak BSD conjecture would be proved for all the elliptic curves  $E_n$ . Due to work of Monsky it is known that every *prime number*  $p \equiv 5, 7 \pmod{8}$  is a congruent number.



# 5 Partitions

## 5.1 The partition function

A partition of a natural number  $n$  is an expression of  $n$  as a sum of natural numbers. More precisely, we define:

**Definition 5.1.1.** A *partition* of a natural number  $n$  is a non-increasing finite sequence of natural numbers  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k > 0$  such that

$$n = \lambda_1 + \lambda_2 + \dots + \lambda_k.$$

The numbers  $\lambda_1, \dots, \lambda_k$  are called the *parts* of the partition.

For example, the partitions of 4 are

$$\begin{aligned} 4 &= 3 + 1 \\ &= 2 + 2 \\ &= 2 + 1 + 1 \\ &= 1 + 1 + 1 + 1. \end{aligned}$$

We will also count the “trivial” partition  $4 = 4$ , so 4 has 5 different partitions. Also note that we do not care about the order of the parts, e.g.  $2 + 1 + 1$  and  $1 + 2 + 1$  are only counted once.

**Definition 5.1.2.** The *partition function*  $p(n)$  is given by the number of partitions of  $n$ .

For convenience, we define  $p(0) = 1$ . Here is a small table with a few values:

$n$	0	1	2	3	4	5	6	...	100	...	1000
$p(n)$	1	1	2	3	5	7	11	...	190 569 292	...	24 061 467 864 032 622 473 692 149 727 991

The partition function  $p(n)$  grows very fast if  $n$  gets large, and listing all partitions becomes unfeasible. One way to compute  $p(n)$  efficiently is by using recursions for  $p(n)$ , which we will derive in this section.

In order to study  $p(n)$ , it is useful to consider some variants which count certain *restricted partitions* of  $n$ :

- $p_d(n)$ : partitions of  $n$  into distinct parts.
- $p_d^{\text{even}}(n)$  and  $p_d^{\text{odd}}(n)$ : partitions of  $n$  into an even or odd number of distinct parts, respectively.
- $p(n, k)$ : partitions of  $n$  into exactly  $k$  parts.

## 5 Partitions

A way to visualize partitions is the *Ferrers diagram*. For example, the partition  $22 = 8 + 4 + 3 + 3 + 2 + 1 + 1$  has the following Ferrers diagram.



Figure 5.1: The Ferrers diagram of  $22 = 8 + 4 + 3 + 3 + 2 + 1 + 1$ .

By reading the Ferrers diagram row-wise, we obtain the corresponding partition of  $n$ . However, we can also read the Ferrers diagram column-wise, which also gives a partition of  $n$ , called the *conjugate* partition. For example, the conjugate partition of  $22 = 8 + 4 + 3 + 3 + 2 + 1 + 1$  is  $22 = 7 + 5 + 4 + 2 + 1 + 1 + 1 + 1$ , corresponding to the “flipped” Ferrers diagram:

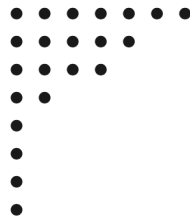


Figure 5.2: The conjugate Ferrers diagram of  $22 = 8 + 4 + 3 + 3 + 2 + 1 + 1$ .

As an application of this idea, we prove the following relation:

**Lemma 5.1.3.** *The function  $p(n, k)$  also counts the number of partitions of  $n$  whose largest part equals  $k$ .*

*Proof.* Let  $P(n, k)$  be the set of partitions of  $n$  into exactly  $k$  parts, such that  $p(n, k) = |P(n, k)|$ , and let  $S(n, k)$  be the set of partitions of  $n$  whose largest part equals  $k$ . We give a bijection between  $P(n, k)$  and  $S(n, k)$ . Consider a partition in  $P(n, k)$ , and draw its Ferrers diagram. It has precisely  $k$  rows, which means that the conjugate partition has largest part equal to  $k$ . Hence the conjugate partition is in  $S(n, k)$ , so conjugation gives a bijection from  $P(n, k)$  to  $S(n, k)$ , and the two sets have equal size.  $\square$

By grouping the partitions of  $n$  by their number of parts, we obtain

$$p(n) = \sum_{k=1}^n p(n, k).$$

We have the following recursion for  $p(n, k)$ , which allows for the quick computation of  $p(n)$ .

**Theorem 5.1.4** (Euler). *We have*

$$p(n, k) = p(n - 1, k - 1) + p(n - k, k).$$

## 5.2 Generating functions and Euler's Pentagonal Theorem

*Proof.* We have seen in the last lemma that  $p(n, k)$  also counts to the number of partitions of  $n$  whose largest part is equal to  $k$ . We let  $S(n, k)$  be the set of all partitions of  $n$  whose largest part is equal to  $k$ , such that  $|S(n, k)| = p(n, k)$ . We will give a bijection

$$S(n, k) \xrightarrow{\sim} S(n-1, k-1) \cup S(n-k, k).$$

Note that the union on the right is disjoint, so this bijection will imply the result. Let  $(\lambda_1, \lambda_2, \dots, \lambda_m) \in S(n, k)$  be a partition of  $n$  with largest part  $\lambda_1 = k$ . If  $\lambda_2 < k$ , then we associate to it the partition

$$(k-1, \lambda_2, \dots, \lambda_m) \in S(n-1, k-1),$$

and if  $\lambda_2 = k$ , then we associate to it the partition

$$(\lambda_2, \lambda_3, \dots, \lambda_m) \in S(n-k, k).$$

This gives the desired bijection. □

## 5.2 Generating functions and Euler's Pentagonal Theorem

In this section we prove Euler's Pentagonal Theorem and use it to derive another recursion for  $p(n)$ . To this end, we consider the generating function of  $p(n)$ . The *generating function* of a sequence  $a_0, a_1, a_2, \dots$  of complex numbers is defined by

$$G(x) = \sum_{n=0}^{\infty} a_n x^n.$$

Here,  $x$  is a variable, and we will not care about the convergence of the series. We will rather use generating functions as a bookkeeping device to facilitate computations with recursively defined sequences.

The generating function of the partition function is given as follows.

**Proposition 5.2.1.** *We have*

$$\sum_{n=0}^{\infty} p(n)x^n = \prod_{n=1}^{\infty} \frac{1}{1-x^n}.$$

*Proof.* Using the geometric series, we can write

$$\begin{aligned} \prod_{n=1}^{\infty} \frac{1}{1-x^n} &= \prod_{n=1}^{\infty} \left( \sum_{k=0}^{\infty} x^{kn} \right) \\ &= (1 + x^{1 \cdot 1} + x^{2 \cdot 1} + x^{3 \cdot 1} + \dots)(1 + x^{1 \cdot 2} + x^{2 \cdot 2} + x^{3 \cdot 2} \dots)(1 + x^{1 \cdot 3} + x^{2 \cdot 3} + x^{3 \cdot 3} + \dots) \dots \end{aligned}$$

Multiplying out, we see that we get a contribution  $+1$  to  $x^n$  from each product of monomials of the form

$$x^{k_1 \cdot 1} x^{k_2 \cdot 2} \dots x^{k_n \cdot n} \quad \text{with} \quad k_1 \cdot 1 + k_2 \cdot 2 + \dots + k_n \cdot n = n.$$

Each such tuple  $(k_1, \dots, k_n)$  corresponds to a unique partition of the form

$$n = \underbrace{1 + \dots + 1}_{k_1 \text{ times}} + \underbrace{2 + \dots + 2}_{k_2 \text{ times}} + \underbrace{3 + \dots + 3}_{k_3 \text{ times}} + \dots$$

Hence, the coefficient at  $x^n$  is  $p(n)$ . □

## 5 Partitions

As a first application of the above generating function identity, we obtain another recursion for  $p(n)$ .

**Proposition 5.2.2.** *We have the recursion*

$$p(n) = \frac{1}{n} \sum_{k=1}^n \sigma(k)p(n-k),$$

where  $\sigma(k) = \sum_{d|k} d$  is the sum of divisors of  $k$ .

*Proof.* By taking the derivative on both sides of the generating function identity for  $p(n)$ , we obtain after a short computation that

$$\sum_{n=1}^{\infty} np(n)x^n = \left( \sum_{n=0}^{\infty} p(n)x^n \right) \cdot \left( \sum_{n=1}^{\infty} \frac{nx^n}{1-x^n} \right).$$

It is easy to check that we have

$$\sum_{n=1}^{\infty} \frac{nx^n}{1-x^n} = \sum_{n=1}^{\infty} \sigma(n)x^n,$$

so we get

$$\sum_{n=1}^{\infty} np(n)x^n = \left( \sum_{n=0}^{\infty} p(n)x^n \right) \cdot \left( \sum_{n=1}^{\infty} \sigma(n)x^n \right).$$

Multiplying out on the right gives

$$\sum_{n=1}^{\infty} np(n)x^n = \sum_{n=1}^{\infty} \left( \sum_{k=1}^n p(n-k)\sigma(k) \right) x^n.$$

Taking the  $n$ -th coefficient on both sides, we obtain the stated recursion for  $p(n)$ . □

Similarly as above, one finds the generating functions

$$\sum_{n=0}^{\infty} p_d(n)x^n = \prod_{n=1}^{\infty} (1+x^n)$$

and

$$1 + \sum_{n=1}^{\infty} (p_d^{\text{even}}(n) - p_d^{\text{odd}}(n))x^n = \prod_{n=1}^{\infty} (1-x^n).$$

They can be used to prove the following important result.

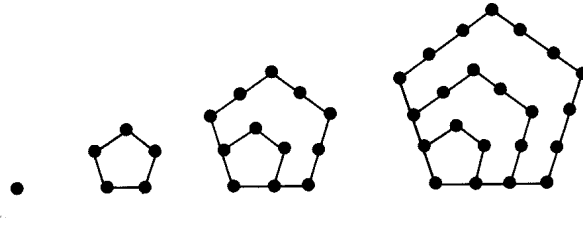
**Theorem 5.2.3** (Euler's Pentagonal Theorem). *We have*

$$\prod_{n=1}^{\infty} (1-x^n) = \sum_{n=-\infty}^{\infty} (-1)^n x^{n(3n-1)/2}.$$

**Remark 5.2.4.** The exponents  $T_n = n(3n-1)/2$  for  $n \in \mathbb{Z}$  are called the (*generalized*) *pentagonal numbers*. For  $n \in \mathbb{N}$  the pentagonal numbers  $T_n = 1, 5, 12, 22, \dots$  count the distinct dots on the edges of  $n-1$  nested pentagons:



## 5.2 Generating functions and Euler's Pentagonal Theorem



*Proof of Euler's Pentagonal Theorem.* Using the generating function

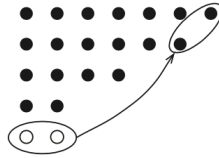
$$\prod_{n=1}^{\infty} (1 - x^n) = 1 + \sum_{n=1}^{\infty} (p_d^{\text{even}}(n) - p_d^{\text{odd}}(n))x^n,$$

we see that the theorem is equivalent to the formula

$$p_d^{\text{even}}(n) - p_d^{\text{odd}}(n) = \begin{cases} (-1)^k, & \text{if } n = k(3k \pm 1)/2, \\ 0, & \text{otherwise.} \end{cases}$$

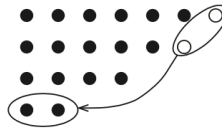
We will give a combinatorial proof of this identity. Using Ferrers diagrams we will construct a map between the sets of partitions of  $n$  into an even or odd number of distinct parts. Consider the Ferrers diagram of a partition of  $n$ , with  $b$  dots in the last row and  $k$  dots on the right-most NE-SW diagonal. We distinguish four cases:

- If  $b < k$ , remove the  $b$  dots on the bottom row and add them as the new right-most diagonal.



This process transforms a partition with an even number of distinct parts into a partition with an odd number of distinct parts and vice versa.

- If  $b > k + 1$ , remove the right-most diagonal and add it as the new bottom row.



Again, this transforms a partition with an even number of distinct parts into a partition with an odd number of distinct parts and vice versa.

- If  $b = k$ , then there is one Ferrers diagram where the above process does not work, of the shape



## 5 Partitions

In this case, we have

$$n = k + (k + 1) + \dots + (2k - 1) = k(3k - 1)/2.$$

- If  $b = k + 1$ , then there is one Ferrers diagram where the above process does not work,



In this case, we have

$$n = (k + 1) + (k + 2) + \dots + 2k = k(3k + 1)/2.$$

If  $n$  is not of the form  $k(3k \pm 1)/2$ , then each Ferrers diagram either has  $b < k$  or  $b > k + 1$ , and the first two transformations give a bijection between the partitions of  $n$  into an even and odd number of distinct parts, hence

$$p_d^{\text{even}}(n) - p_d^{\text{odd}}(n) = 0.$$

If  $n = k(3k \pm 1)/2$ , then  $b$  and  $k$  are determined from  $n$ , so there is precisely one exceptional Ferrers diagram which cannot be transformed with the first two steps. This diagram has  $k$  rows, so the partition has  $k$  parts, which gives

$$p_d^{\text{even}}(n) - p_d^{\text{odd}}(n) = (-1)^k.$$

This finishes the proof. □

As an application, we obtain another recursion for the partition number  $p(n)$ .

**Theorem 5.2.5** (Euler). *We have*

$$\begin{aligned} p(n) &= \sum_{k=1}^{\infty} (-1)^{k+1} \left[ p\left(n - \frac{3k^2 - k}{2}\right) + p\left(n - \frac{3k^2 + k}{2}\right) \right] \\ &= p(n - 1) + p(n - 2) - p(n - 5) - p(n - 7) + p(n - 12) + p(n - 15) - \dots \end{aligned}$$

*Proof.* Using the generating functions

$$\prod_{n=1}^{\infty} \frac{1}{1 - x^n} = \sum_{\ell=0}^{\infty} p(\ell) x^{\ell}$$

and

$$\prod_{n=1}^{\infty} (1 - x^n) = 1 + \sum_{k=1}^{\infty} (-1)^k \left( x^{k(3k+1)/2} + x^{k(3k-1)/2} \right),$$

we obtain

$$1 = \left( \sum_{\ell=0}^{\infty} p(\ell) x^{\ell} \right) \cdot \left( 1 + \sum_{k=1}^{\infty} (-1)^k \left( x^{k(3k+1)/2} + x^{k(3k-1)/2} \right) \right).$$

Hence, the  $n$ -th coefficient for  $n > 0$  of the right-hand side vanishes. On the other-hand, it is given by

$$p(n) + \sum_{k=1}^{\infty} (-1)^k \left[ p\left(n - \frac{k(3k+1)}{2}\right) + p\left(n - \frac{k(3k-1)}{2}\right) \right].$$

This gives the stated formula. □