

## Chapter 5

# Pell's Equation

One of the earliest issues grappled with in number theory is the fact that geometric quantities are often not rational. For instance, if we take a right triangle with two side lengths equal to 1, the hypotenuse has length  $\sqrt{2}$ , which is irrational. But how can we do arithmetic with irrational numbers? Well, perhaps the most basic thing is to work with rational approximations. Almost 4000 years ago, Babylonians had discovered the following approximation to  $\sqrt{2}$ :

$$\sqrt{2} = 1.41421356\dots \approx \frac{30547}{21600} = 1.41421\overline{296}. \quad (5.0.1)$$

In this chapter we'll explain how to find the (integer) solutions to **Pell's equation**:

$$x^2 - dy^2 = 1, \quad (5.0.2)$$

and how this gives us good approximations to  $\sqrt{d}$  (see [Proposition 5.2.3](#)).

Following Stigler's law of eponymy<sup>1</sup>, Pell's equation was studied by the Indian mathematician and astronomer Brahmagupta in 628 (who discovered the composition law [Proposition 5.1.1](#)) and with a general method of solution by another Indian mathematician and astronomer, Bhaskara II, in 1150. In Europe, methods for solving Pell's equation were rediscovered hundreds of years later by Fermat and Lord Brouncker. Euler misattributed Lord Brouncker's solution after reading a discussion of Lord Brouncker's method written by the English mathematician John Pell (1611–1685).

Throughout this chapter  $d > 1$  is a positive integer which is not a square.

### 5.1 Units and Pell's equation

Recall that a unit  $u$  of  $\mathbb{Z}[\sqrt{d}]$  was defined to be an element such that  $u$  has a multiplicative inverse  $u^{-1} \in \mathbb{Z}[\sqrt{d}]$  (i.e., the real number  $u \neq 0$  and  $\frac{1}{u} \in \mathbb{Z}[\sqrt{d}]$ ). Further, by [Lemma 2.1.2](#), we know that  $x + y\sqrt{d} \in \mathbb{Z}[\sqrt{d}]$  ( $x, y \in \mathbb{Z}$ ) is a unit if and only if

$$N(x + y\sqrt{d}) = (x + y\sqrt{d})(x - y\sqrt{d}) = x^2 - dy^2 = \pm 1.$$

---

<sup>1</sup>That no scientific discovery is named after its first discoverer. The Pythagorean theorem is another famous example. Of course there are many counterexamples to Stigler's law as well. Appropriately, Stigler's law itself is not.

Thus solutions to Pell's equation (5.0.2) are in natural bijection with the units of  $\mathbb{Z}[\sqrt{d}]$  with norm 1.

On the other hand, we also know by Proposition 3.3.4 that the set of units

$$U = U_d = \mathbb{Z}[\sqrt{d}]^\times$$

of  $\mathbb{Z}[\sqrt{d}]$  form an (abelian) group. We also denote by  $U^+ = U_d^+$  the set of units in  $U = U_d$  of norm 1, so we can think of the solutions to Pell's equation as the subgroup  $U^+$  of  $U$ .

**Exercise 5.1.1.** Check that  $U^+$  is indeed a subgroup of  $U$ .

This group structure will help us determine the set of solutions to Pell's equation. First, we have the following, which is similar to the composition law for sums of two squares.

**Proposition 5.1.1. (Composition law)** *If  $(x_1, y_1)$  and  $(x_2, y_2)$  are solutions to*

$$x_1^2 - dy_1^2 = m, \quad x_2^2 - dy_2^2 = n.$$

*Then the composition of these solutions defined by*

$$(x_3, y_3) = (x_1, y_1) \cdot (x_2, y_2) := (x_1x_2 + dy_1y_2, x_1y_2 + y_1x_2) \quad (5.1.1)$$

*is a solution of*

$$x_3^2 - dy_3^2 = mn.$$

*Proof.* We simply translate the above into a statement about norms. The hypothesis says  $N(x_1 + y_1\sqrt{d}) = m$  and  $N(x_2 + y_2\sqrt{d}) = n$ . Now observe that

$$(x_1 + y_1\sqrt{d})(x_2 + y_2\sqrt{d}) = x_1x_2 + ny_1y_2 + (x_1y_2 + y_1x_2)\sqrt{d} = x_3 + y_3\sqrt{d}.$$

Hence by the multiplicative property of the norm,

$$x_3^2 - dy_3^2 = N(x_3 + y_3\sqrt{d}) = N(x_1 + y_1\sqrt{d})N(x_2 + y_2\sqrt{d}) = mn.$$

□

In particular, when  $m = n = 1$ , this says that we can compose two solutions to Pell's equation to get a third solutions. We can also compose two solutions to  $x^2 - dy^2 = -1$  to get a solution to  $x^2 - dy^2 = +1$ . Both of these are summarized in this corollary.

**Corollary 5.1.2.** *Let  $u_1 = x_1 + y_1\sqrt{d}$  and  $u_2 = x_2 + y_2\sqrt{d}$  be units of  $\mathbb{Z}[\sqrt{d}]$ . If  $N(u_1) = N(u_2)$ , then the composition  $(x_1, y_1) \cdot (x_2, y_2)$  defined in Eq. (5.1.1) also a solution to Pell's equation (5.0.2).*

The following should be obvious if you've had an algebra class, but since we never covered isomorphisms, I'm not entirely sure if this is obvious to you:

**Exercise 5.1.2.** Let  $G \subset \mathbb{Z} \times \mathbb{Z}$  be the set of solutions to Pell's equation (5.0.2). Show that the composition law Eq. (5.1.1) makes  $G$  into a group. (Note the above corollary just says composition is a binary operation on  $G$ .)

Now that we know the composition law, the hope is that if we can determine a few good solutions to Pell's equation, then maybe we can generate all solutions by composing those we know. Moreover, ideally these good solutions should be the “smallest nontrivial” solutions to Pell's equation. By the **trivial solutions** to Pell's equation, we mean the obvious ones:  $(x, y) = (\pm 1, 0)$ , which correspond to the elements of  $U^+$  which lie in  $\mathbb{Z}$ , i.e.,  $\pm 1$ .

**Example 5.1.1.** Consider  $d = 2$ . Note  $(1, 1)$  is a solution to  $x^2 - 2y^2 = -1$ . The composition  $(1, 1) \cdot (1, 1) = (3, 2)$  is a nontrivial solution to Pell's equation  $x^2 - 2y^2 = 1$ . Similarly, we compute  $(3, 2) \cdot (3, 2) = (17, 12)$  and  $(3, 2) \cdot (17, 12) = (99, 70)$ .

Hence  $(3, 2)$ ,  $(17, 12)$  and  $(99, 70)$  are three nontrivial solutions to  $x^2 - 2y^2 = 1$ .

**Exercise 5.1.3.** Find a nontrivial solution to  $x^2 - 3y^2 = 1$ . Use composition to find two more (distinct) solutions to  $x^2 - 3y^2 = 1$ .

We remark that in the case of  $d = 3$ , unlike  $d = 2$ , there are no units of norm  $-1$ . In fact, the following more general statement is true.

**Exercise 5.1.4.** Suppose  $d \equiv 3 \pmod{4}$ . Show  $\mathbb{Z}[\sqrt{d}]$  has no units of norm  $-1$ , i.e.,  $U_d = U_d^+$ .

The converse to the previous exercise does not hold, i.e., there may or may not be a unit of norm  $-1$  when  $d \not\equiv 3 \pmod{4}$ . We've seen there is such a unit when  $d = 2$ . The next exercise gives you an example where there isn't a unit of norm  $-1$  but  $d \not\equiv 3 \pmod{4}$ .

**Exercise 5.1.5.** Show that  $\mathbb{Z}[\sqrt{6}]$  has no units of norm  $-1$ .

In general, it is an open problem to determine for what  $d$  there are units of norm  $-1$  in  $\mathbb{Z}[\sqrt{d}]$ . It's not clear that there is a nice answer to this problem, but there results about how often  $\mathbb{Z}[\sqrt{d}]$  has units of norm  $-1$ .

## 5.2 Approximation and existence of solutions

At the end of the last section we saw that there are nontrivial solutions to Pell's equation when  $d = 2, 3$ . Next we will prove the existence of a non-trivial solution for all nonsquare  $d$ , which is originally due to Lagrange in 1768. However, the proof we will give is due to Dirichlet (ca. 1840). It uses the pigeonhole principle. You've probably at least seen the finite version in your Discrete Math class.

**Pigeonhole principle**

- (finite version) If  $m > k$  pigeons go into  $k$  boxes, at least one must box must contain more than 1 pigeon.
- (infinite version) If infinitely many pigeons go into  $k$  boxes, at least one box must contain infinitely many pigeons).

**Proposition 5.2.1. (Dirichlet's approximation theorem)** *For any nonsquare  $d > 1$  and integer  $B > 1$ , there exist  $a, b \in \mathbb{N}$  such that  $b < B$  and*

$$|a - b\sqrt{d}| < \frac{1}{B}.$$

This says that  $|\frac{a}{b} - \sqrt{d}| < \frac{1}{bB}$ , which is a precise way of saying  $\frac{a}{b}$  is close to  $\sqrt{d}$ . E.g., it says we can find a rational approximation  $\frac{a}{b}$  for  $\sqrt{2}$  which is accurate within  $\frac{1}{100,000}$  (so  $\frac{a}{b}$  and  $\sqrt{2}$  agree to 5 decimal places, after rounding if necessary<sup>2</sup>) with denominator  $b < 100,000$ . Such an example was exhibited at the beginning of this chapter in (5.0.1).

*Proof.* Consider the  $B - 1$  irrational numbers

$$\sqrt{d}, 2\sqrt{d}, \dots, (B - 1)\sqrt{d}.$$

For each such number  $k\sqrt{d}$  ( $1 \leq k \leq B - 1$ ), let  $a_k \in \mathbb{N}$  be such that

$$0 < a_k - k\sqrt{d} < 1.$$

Partition the interval  $[0, 1]$  into  $B$  subintervals of length  $\frac{1}{B}$ . Then, of the  $B + 1$  numbers

$$0, a_1 - \sqrt{d}, a_2 - \sqrt{d}, \dots, a_{B-1} - (B - 1)\sqrt{d}, 1$$

in  $[0, 1]$  two of them must be in the same subinterval of length  $\frac{1}{B}$ . Hence they are less than distance  $\frac{1}{B}$  apart, i.e., their difference satisfies  $|a - b\sqrt{d}| < \frac{1}{B}$ . Further their irrational parts must be distinct, so we have  $-B < b < B$  with  $b \neq 0$ . If  $b > 0$  we are done; if  $b < 0$ , simply multiply  $a$  and  $b$  by  $-1$ . Clearly we need  $a > 0$  for  $|a - b\sqrt{d}| < 1$ .  $\square$

**Theorem 5.2.2.** *Suppose  $d \in \mathbb{N}$  is nonsquare. Then  $x^2 - dy^2 = 1$  has a nontrivial solution, i.e., there is a unit in  $\mathbb{Z}[\sqrt{d}]$  of norm 1 other than  $\pm 1$ .*

*Proof. Step 1.* Fix  $B_1 > 1$ . Then by Dirichlet's approximation theorem, there exist  $a_1, b_1 \in \mathbb{N}$  such that  $|a_1 - b_1\sqrt{d}| < \frac{1}{B_1} < \frac{1}{b_1}$ . Let  $B_2 > b_1$  such that  $\frac{1}{B_2} < |a_1 - b_1\sqrt{d}|$ . Applying Dirichlet's approximation again, we get a new pair  $(a_2, b_2)$  of integers such that

$$|a_2 + b_2\sqrt{d}| < \frac{1}{B_2} < \frac{1}{b_2}.$$

<sup>2</sup>For instance, 0.5006 and 0.49998 are within  $\frac{1}{1000}$  of each other, but their first three digits are only equal after rounding to the nearest 4 digits.

Repeating this we see there an infinite sequence of distinct integer pairs  $(a_j, b_j)$  such that  $|a_j - b_j\sqrt{d}|$  gets smaller and smaller, and

$$|a_j - b_j\sqrt{d}| < \frac{1}{b_j}.$$

for all  $j \geq 1$ . Then  $\frac{a_j}{b_j}$  is an infinite sequence of increasingly good approximations to  $\sqrt{d}$ .

**Step 2.** Assume  $(a, b)$  satisfies  $|a - b\sqrt{d}| < \frac{1}{b}$ . Note that

$$|a + b\sqrt{d}| \leq |a - b\sqrt{d}| + |2b\sqrt{d}| \leq 1 + 2b\sqrt{d} \leq 3b\sqrt{d}.$$

Then

$$|a^2 - db^2| = |a + b\sqrt{d}||a - b\sqrt{d}| \leq 3b\sqrt{d}\frac{1}{b} = 3\sqrt{d}.$$

Hence there are infinitely many  $a - b\sqrt{d} \in \mathbb{Z}[\sqrt{d}]$  whose norm, in absolute values, is at most  $3\sqrt{d}$ .

**Step 3.** By successive applications of the (infinite) pigeonhole principle, we have:

- (i) infinitely many  $a - b\sqrt{d}$  with the same norm  $n \in \mathbb{Z}$ , where  $|n| \leq 3\sqrt{d}$  (and  $n \neq 0$ )
- (ii) infinitely many  $a - b\sqrt{d}$  with norm  $n$  and  $a \equiv a_0 \pmod{n}$  for some  $a_0$ .
- (iii) infinitely many  $a - b\sqrt{d}$  with norm  $n$ ,  $a \equiv a_0 \pmod{n}$ ,  $b \equiv b_0 \pmod{n}$  for some  $b_0$ .

Hence, relabeling if necessary, we have  $a_1, b_1, a_2, b_2 \in \mathbb{N}$  such that  $N(a_1 - b_1\sqrt{d}) = N(a_2 - b_2\sqrt{d}) = n$ ,  $a_1 \equiv a_2 \pmod{n}$ ,  $b_1 \equiv b_2 \pmod{n}$ , and  $a_1 - b_1\sqrt{d} \neq \pm(a_2 - b_2\sqrt{d})$ .

**Step 4.** Consider

$$\alpha := \frac{a_1 - b_1\sqrt{d}}{a_2 - b_2\sqrt{d}} = \frac{(a_1 - b_1\sqrt{d})(a_2 + b_2\sqrt{d})}{a_2^2 - db_2^2} = \frac{a_1a_2 - db_1b_2}{n} + \frac{a_1b_2 - b_1a_2}{n}\sqrt{d}.$$

Note

$$a_1a_2 - db_1b_2 \equiv a_1a_1 - db_1b_1 \equiv a_1^2 - db_1^2 \equiv 0 \pmod{n},$$

and

$$a_1b_2 - b_1a_2 \equiv a_1b_1 - b_1a_1 \equiv 0 \pmod{n}.$$

Thus the coefficients of  $\alpha$  are integers, i.e.,  $\alpha = a + b\sqrt{d} \in \mathbb{Z}[\sqrt{d}]$  where  $a, b \in \mathbb{Z}$ . Then hence since

$$a^2 - db^2 = N(a + b\sqrt{d}) = N(a_1 - b_1\sqrt{d})N\left((a_2 - b_2\sqrt{d})^{-1}\right) = nn^{-1} = 1,$$

i.e.,  $(a, b)$  is a solution of  $x^2 - dy^2 = 1$ . Furthermore, it is a nontrivial solution since  $a_1 - b_1\sqrt{d} \neq \pm(a_2 - b_2\sqrt{d})$ .  $\square$

**Exercise 5.2.1.** Explain how to modify the above proof to conclude the existence of infinitely many solutions to  $x^2 - dy^2 = 1$ . Conclude the real quadratic rings  $\mathbb{Z}[\sqrt{d}]$  ( $d > 1$  nonsquare) have infinitely many units, in contrast to the case of imaginary quadratic rings  $\mathbb{Z}[\sqrt{-d}]$ .

The above approximations suggest how solutions to Pell's equation are related to rational approximations to  $\sqrt{d}$ .

**Proposition 5.2.3.** *Suppose  $(x, y)$  is a nontrivial solution to  $x^2 - dy^2 = 1$ . Assume  $x, y > 0$ . Then*

$$0 < \frac{x}{y} - \sqrt{d} < \frac{1}{y(1 + \sqrt{d})} < \frac{1}{2y}.$$

*Proof.* Let  $\alpha = x - y\sqrt{d} \in \mathbb{Z}[\sqrt{d}]$ . Then  $N(\alpha) = \alpha\bar{\alpha} = 1$ . Note  $\bar{\alpha} = x + y\sqrt{d} \geq 1 + \sqrt{d}$  since  $x, y \geq 1$ . Thus  $\alpha \leq \frac{1}{1 + \sqrt{d}}$ , so

$$\frac{x}{y} - \sqrt{d} = \frac{\alpha}{y} < \frac{1}{y(1 + \sqrt{d})}.$$

Since  $\bar{\alpha}$  and  $N(\alpha)$  are positive, so is  $\alpha$ , and thus  $\frac{\alpha}{y}$ , which finishes the asserted bounds.  $\square$

This proposition says positive solutions  $(x, y)$  to Pell's equation, i.e., units of norm +1, give rational approximations to  $\sqrt{d}$ , and solutions with larger values of  $y$  give better approximations. Furthermore, we are always getting overestimates for  $\sqrt{d}$ . One can similarly get underestimates with units of norm  $-1$ , i.e., solutions to  $x^2 - dy^2 = -1$ , at least when they exist. When they don't, one could instead look for solutions to  $x^2 - dy^2 = -2$  or  $x^2 - dy^2 = -3$  etc. We remark the approximation in (5.0.1) corresponds to the solution  $(30547, 21600)$  to  $x^2 - 2y^2 = -791$ , though I'm not suggesting Babylonians came up with this approximation by starting with the equation  $x^2 - 2y^2 = -791$ !

**Exercise 5.2.2.** Suppose  $(x, y)$  is solution to  $x^2 - dy^2 = -1$  with  $x, y > 0$ . Show  $\frac{x}{y} < \sqrt{d}$  and prove a (good) bound for  $\sqrt{d} - \frac{x}{y}$  in terms of  $y$ .

**Exercise 5.2.3.** Suppose  $(x, y)$  is solution to  $x^2 - dy^2 = -2$  with  $x, y > 0$ . Show  $\frac{x}{y} < \sqrt{d}$  and prove a (good) bound for  $\sqrt{d} - \frac{x}{y}$  in terms of  $y$ .

**Example 5.2.1.** Recall from Example 5.1.1,  $(3, 2)$ ,  $(17, 12)$  and  $(99, 70)$  are solutions to  $x^2 - 2y^2 = 1$ . This gives the following approximations, with the following error bounds  $\frac{1}{2y}$  from the above proposition:

$\frac{x}{y}$	decimal	error bound	$\frac{x}{y} - \sqrt{2}$
$\frac{3}{2}$	1.5	$< 0.25$	0.085786...
$\frac{17}{12}$	1.41 $\bar{6}$	$< 0.041\bar{6}$	0.0024531...
$\frac{99}{70}$	1.4142857	$< 0.0071428\bar{5}$	0.00007215...

**Exercise 5.2.4.** Recall the solutions in the previous example came from the first three powers  $\epsilon^n$  ( $n = 1, 2, 3$ ) of the unit  $\epsilon = 3 - 2\sqrt{2} \in U_2^+$ . However, none of these approximations are better than the (more complicated) Babylonian one (5.0.1). Using a calculator, compute a few successive approximations to  $\sqrt{2}$  from higher powers  $\epsilon^n$ , along with the exact error (up to several decimal places). What is the first approximation you get in this way that is better (i.e., closer to  $\sqrt{2}$ ) than the one in (5.0.1).

**Exercise 5.2.5.** Using 3 nontrivial solutions to  $x^2 - 3y^2 = 1$  you found in Exercise 5.1.3, give 3 rational approximations to  $\sqrt{3}$  with error bounds. Using a calculator, compute the actual error in these approximations (e.g., you can make a table as in the example above).

### 5.3 Fundamental units

Here we determine the structure of the group of units  $U_d$ , which will give us a method for generating all solutions to Pell's equation.

**Definition 5.3.1.** The **fundamental unit**  $\epsilon_d$  of  $\mathbb{Z}[\sqrt{d}]$  is the smallest unit  $x + y\sqrt{d} \in \mathbb{Z}[\sqrt{d}]$  such that  $x, y > 0$ . The **fundamental +unit**  $\epsilon_d^+$  of  $\mathbb{Z}[\sqrt{d}]$  is the smallest unit  $x + y\sqrt{d} \in \mathbb{Z}[\sqrt{d}]$  such that  $x, y > 0$  and  $N(\epsilon) = 1$ .<sup>3</sup>

In real quadratic rings, smallest means with respect to the usual order on  $\mathbb{R}$ , unlike the case of imaginary quadratic rings where we (partially) ordered elements by their norm.

**Lemma 5.3.2.** For any (nonsquare)  $d > 1$ , the fundamental unit  $\epsilon_d$  and the fundamental +unit  $\epsilon_d^+$  exist and are uniquely defined.

*Proof.* Since  $<$  defines a strict ordering of real numbers, the condition of “smallest” guarantees that  $\epsilon_d$  and  $\epsilon_d^+$  will be unique if they exist, so it suffices to show existence.

Recall we always have a nontrivial solution  $(x_0, y_0)$  to  $|x^2 - dy^2| = 1$  from Theorem 5.2.2. Moreover, we can assume  $x_0, y_0 > 0$ . Now note if  $x + y\sqrt{d} < \epsilon_0 = x_0 + y_0\sqrt{d}$  with  $x, y > 0$ , we must have  $x < \epsilon_0$  and  $y < \frac{\epsilon_0}{\sqrt{d}}$ . Hence

$$\epsilon_d = \min \left\{ x + y\sqrt{d} : 1 \leq x, \sqrt{d}y < \epsilon_0, |x^2 + dy^2| = 1 \right\}.$$

Since the set on the right is finite, this minimum is well defined, hence  $\epsilon_d$  exists.

The case of  $\epsilon_d^+$  follows in the same way, simply using the equation  $x^2 - dy^2 = 1$  instead of  $|x^2 - dy^2| = 1$ .  $\square$

<sup>3</sup>Note that most discussions you will find about fundamental units talk about fundamental units in the ring of integers  $\mathcal{O}_d$  of  $\mathbb{Q}(\sqrt{d})$ . Here, assuming  $d$  is squarefree,  $\mathcal{O}_d$  is just  $\mathbb{Z}[\sqrt{d}]$  when  $d \not\equiv 1 \pmod{4}$  but is  $\mathbb{Z}[\frac{1+\sqrt{d}}{2}]$  when  $d \equiv 1 \pmod{4}$  (recall Definition 2.5.5). So be careful comparing what we say here and what is written other places about fundamental units, as there may be a slight difference when  $d \equiv 1 \pmod{4}$ , though there is no serious difference in the theory. E.g., when  $d = 5$  a fundamental unit in  $\mathbb{Z}[\sqrt{5}]$  is  $2 + \sqrt{5}$  but in  $\mathcal{O}_5$  it is  $\frac{1+\sqrt{5}}{2}$ . On the other hand, the fundamental unit in  $\mathcal{O}_{17} = \mathbb{Z}[\frac{1+\sqrt{17}}{2}]$  is the same as the fundamental unit in  $\mathbb{Z}[\sqrt{17}]$ , namely  $\epsilon_{17} = 4 + \sqrt{17}$ .

Also, the term “fundamental +unit” is not standard—as far as I know, there is no standard term for the phrase “the smallest unit of norm 1.”

Here is a naive algorithm for finding  $\varepsilon_d$  or  $\varepsilon_d^+$ . First, pick some bound  $N \geq 1$ . Then range over all  $1 \leq x, y \leq N$  to look for solutions to  $|x^2 - dy^2| = 1$  or  $x^2 - dy^2 = 1$ . If we found any, then the smallest solution  $(x, y)$  gives us  $\varepsilon_d$  or  $\varepsilon_d^+$  as  $x + y\sqrt{d}$ . If not, we pick a larger  $N$  and repeat. This process terminates at some point by the existence of  $\varepsilon_d$  and  $\varepsilon_d^+$ .

**Example 5.3.1.** When  $d = 2$ ,  $\varepsilon_2 = 1 + \sqrt{2}$  and  $\varepsilon_2^+ = \varepsilon_2^2 = 3 + 2\sqrt{2}$ . We saw both of these units in [Example 5.1.1](#).

**Example 5.3.2.** When  $d = 5$ , we compute  $\varepsilon_5 = 2 + \sqrt{5}$  and  $\varepsilon_5^+ = \varepsilon_5^2 = 9 + 4\sqrt{5}$ .

**Example 5.3.3.** Consider  $d = 7$ . Recall from [Exercise 5.1.4](#) that  $\mathbb{Z}[\sqrt{7}]$  has no units of norm  $-1$ . We compute  $\varepsilon_7 = \varepsilon_7^+ = 8 + 3\sqrt{7}$ .

**Exercise 5.3.1.** Compute  $\varepsilon_d$  and  $\varepsilon_d^+$  for  $d = 3, 6, 11$ .

**Exercise 5.3.2.** An alternative definition of fundamental unit (resp.  $+$ -unit) is the smallest  $\varepsilon > 1$  in  $\mathbb{Z}[\sqrt{d}]$  such that  $|N(\varepsilon)| = 1$  (resp.  $N(\varepsilon) = 1$ ). Prove that this is equivalent to the above definition as follows. (*Suggestion:* Show that  $\varepsilon = x + y\sqrt{d} > 1$  a unit implies  $|\bar{\varepsilon}| < 1$ , which implies  $x, y > 0$ .)

**Theorem 5.3.3.** For  $d > 1$  nonsquare,  $U_d$  (resp.  $U_d^+$ ) is the infinite abelian group generated by  $\varepsilon_d$  (resp.  $\varepsilon_d^+$ ) and  $-1$ . Explicitly,

$$U_d = \{\dots, \pm\varepsilon_d^{-2}, \pm\varepsilon_d^{-1}, \pm 1, \pm\varepsilon_d, \pm\varepsilon_d^2, \dots\}$$

and

$$U_d^+ = \{\dots, \pm(\varepsilon_d^+)^{-2}, \pm(\varepsilon_d^+)^{-1}, \pm 1, \pm\varepsilon_d^+, \pm(\varepsilon_d^+)^2, \dots\},$$

and all the elements listed in the sets on the right are distinct, i.e.,  $\varepsilon_d^m = \pm\varepsilon_d^n$  for  $m, n \in \mathbb{Z}$  (resp.  $(\varepsilon_d^+)^m = \pm(\varepsilon_d^+)^n$ ) if and only if  $m = n$  and the plus/minus sign is  $+$ .

*Proof.* We know  $U_d$  and  $U_d^+$  are abelian groups by [Proposition 3.3.4](#) and [Exercise 5.1.1](#). Thus  $U_d$  and  $U_d^+$  must contain all the elements in the sets on the right.

Write  $G$  denote  $U_d$  or  $U_d^+$ , and let  $\varepsilon$  denote  $\varepsilon_d$  or  $\varepsilon_d^+$ , according to whether  $G = U_d$  or  $G = U_d^+$ . Since  $\varepsilon > 1$ , the sequence  $\varepsilon^n$  ( $n \geq 0$ ) is a strictly increasing sequence lying in  $[1, \infty)$ , and  $\varepsilon^{-n}$  ( $n > 0$ ) is a strictly decreasing sequence lying in  $(0, 1)$ . From this one easily sees that all elements in the above sets on the right are distinct.

Finally, we show any  $\alpha \in G \subset \mathbb{Z}[\sqrt{d}]$  is of the form  $\pm\varepsilon^n$  for some  $n \in \mathbb{Z}$ . Suppose there is some  $\alpha$  which is not of this form. By taking the negative and/or inverse if necessary, we may assume  $\alpha > 1$ . Since  $\varepsilon$  is the smallest element of  $G$  larger than 1 ([Exercise 5.3.2](#)) and  $\varepsilon^n \rightarrow \infty$  as  $n \rightarrow \infty$ , there must be some  $n > 0$  such that  $\varepsilon^m < \alpha < \varepsilon^{m+1}$ . But then  $1 < \alpha\varepsilon^{-m} < \varepsilon$  and  $N(\alpha\varepsilon^{-m}) = 1$ , contradicting the minimality of  $\varepsilon$ .  $\square$

Note that for any unit  $u \in U_d$ ,  $N(u) = u\bar{u} = \pm 1$  implies  $u^{-1}$  is either  $\bar{u}$  or  $-\bar{u}$ , according to whether  $N(u) = 1$  or  $N(u) = -1$ . In particular, if  $(\varepsilon_d^+)^n = x_n + y_n\sqrt{d}$ , then  $(\varepsilon_d^+)^{-n} = \overline{(\varepsilon_d^+)^n} = x_n - y_n\sqrt{d}$ .

Hence the above theorem says that once we find  $\varepsilon_d^+$ , we can compute all elements of  $U_d^+$  by computing  $(\varepsilon_d^+)^n = x_n + y_n\sqrt{d}$ ,  $n \geq 1$ . Then the elements of  $U_d^+$  are

$$U_d^+ = \{\pm 1\} \cup \left\{ \pm x_n \pm y_n\sqrt{d} : n \geq 1 \right\}, \quad (5.3.1)$$

where we read the  $\pm$  signs in  $\pm x_n \pm y_n$  independently. (A similar statement is also true for  $U_d$ .) This immediately gives our desired description of solutions to (5.0.2).

**Corollary 5.3.4.** *For  $n \geq 1$ , write  $(\varepsilon_d^+)^n = x_n + y_n\sqrt{d}$  for  $n \geq 1$  (with  $x_n, y_n \in \mathbb{Z}$ ). Then all solutions to Pell's equation  $x^2 - dy^2 = 1$  are the trivial solutions  $(\pm 1, 0)$  and the nontrivial solutions  $(\pm x_n, \pm y_n)$  for  $n \geq 1$ .*

Via Proposition 5.2.3, this gives us the following sequence of approximations

$$\frac{x_n}{y_n} \approx \sqrt{d}$$

of  $\sqrt{d}$ . To prove that these approximations are getting better (at least asymptotically), by this proposition we want to prove the  $y_n$ 's are increasing.

**Exercise 5.3.3.** With  $x_n, y_n$  as above, show the sequences  $(x_n)$  and  $(y_n)$  are strictly increasing sequences for  $n \geq 1$ . Deduce that the sequence  $\frac{x_n}{y_n}$  converges to  $\sqrt{d}$ .

Using the above theorem, we can also relate  $\varepsilon_d$  and  $\varepsilon_d^+$  now.

**Exercise 5.3.4.** For  $d > 1$  a nonsquare, show  $\varepsilon_d^+ = \varepsilon_d^2$  if  $\mathbb{Z}[\sqrt{d}]$  has units of norm  $-1$ , and  $\varepsilon_d^+ = \varepsilon_d$  otherwise. Deduce in particular that  $\varepsilon_d \leq \varepsilon_d^+$ .

Hence if we solve the problem of finding the fundamental unit  $\varepsilon_d$ , we also know the fundamental +unit  $\varepsilon_d^+$ . Since  $\varepsilon_d \leq \varepsilon_d^+$ , even if our goal is to compute  $\varepsilon_d^+$ , it may often be easier algorithmically to look for  $\varepsilon_d$  first, since the  $x$  and  $y$  appearing in the representation  $x + y\sqrt{d}$  can be much smaller.

**Example 5.3.4.** Consider  $d = 29$ . Then by the naive algorithm for finding fundamental units, we can check  $\varepsilon_{29} = 70 + 13\sqrt{29}$ , which has norm  $-1$ . Thus  $\varepsilon_{29}^+ = \varepsilon_{29}^2 = 9801 + 1820\sqrt{29}$ , but this would require many more calculations to find solely by the naive algorithm.

Here's another consequence of the structure theorem for  $U_d$  (or, if you prefer, the previous exercise): if  $\mathbb{Z}[\sqrt{d}]$  has no units of norm  $-1$ , we can prove this algorithmically by computing  $\varepsilon_d$  and checking it has norm 1. For then  $\pm\varepsilon_d^n$  also has norm 1 for all  $n$ , i.e.,  $U_d = U_d^+$ .

## 5.4 Continued fractions

In the last section, we described how to find all solutions to Pell's equation in terms of the fundamental +unit  $\varepsilon_d^+$ . Earlier, we also presented a naive algorithm to compute  $\varepsilon_d$  and  $\varepsilon_d^+$ . The problem is that as  $d$  gets even moderately large, the naive algorithm is not very efficient. This is already suggested by the case of  $d = 29$  in [Example 5.3.4](#). Here is a more impressive example:

**Example 5.4.1.** When  $d = 61$ ,  $\varepsilon_d^+ = 1766319049 + 226153980\sqrt{61}$ , i.e., the smallest positive nontrivial solution to  $x^2 - 61y^2 = 1$  is  $(1766319049, 226153980)$ .

The above example was discovered by Bhaskara II in the 12th century in India, and independently (much later) rediscovered by Fermat in Europe. How can one find such solutions, especially without powerful computing devices? The answer come from an alternative representation of numbers, not as decimals, but as *continued fractions*.

First we explain the continued fraction expansion with an example.

**Example 5.4.2.** Consider  $\frac{a}{b} = \frac{a_1}{b_1} = \frac{13}{5}$ , which we write as a whole number plus a remainder:

$$\frac{a_1}{b_1} = \frac{13}{5} = 2 + \frac{3}{5}.$$

Now we can't exactly repeat this on the remainder, but we can on its *reciprocal*:

$$\frac{a_2}{b_2} := \frac{5}{3} = 1 + \frac{2}{3}.$$

Thus we have

$$\frac{a}{b} = 2 + \frac{1}{3/2} = 2 + \frac{1}{1 + \frac{2}{3}}.$$

Now repeat again with the reciprocal of the remainder in  $\frac{a_2}{b_2}$ :

$$\frac{a_3}{b_3} := \frac{3}{2} = 1 + \frac{1}{2}.$$

If we do this again, we get:

$$\frac{a_4}{b_4} = \frac{2}{1} = 2,$$

a rational with no remainder, and so we stop. This leads to the following expression:

$$\frac{a}{b} = \frac{13}{5} = \boxed{2} + \frac{1}{\boxed{1} + \frac{1}{\boxed{1} + \frac{1}{\boxed{2}}}},$$

which we call the continued fraction expansion of  $\frac{13}{5}$ . Note that because at each stage we are taking reciprocals, we'll see a sequence of 1's going down, and all that matters are the

numbers in the boxes. To simplify notation, we will also write this as

$$[2, 1, 1, 2] = 2 + \frac{1}{1 + \frac{1}{1 + \frac{1}{2}}} = 2 + \frac{1}{1 + \frac{1}{1 + \frac{1}{2}}}$$

**Definition 5.4.1.** Let  $x \in \mathbb{R}$ . The **continued fraction expansion** of  $x$  is the expression  $[q_1, q_2, q_3, \dots] = q_1 + \frac{1}{q_2 + \frac{1}{q_3 + \dots}}$  where  $q_1 \in \mathbb{Z}$  and  $q_j \in \mathbb{Z}_{\geq 0}$  for  $j \geq 2$  are defined as follows:

- $q_1 = \lfloor x \rfloor$  is the greatest integer  $\leq x$ , so  $0 \leq r_1 < 1$  where  $r_1 = x - q_1$ ;
- for  $j \geq 1$ , inductively set

$$q_{j+1} = \begin{cases} \lfloor \frac{1}{r_j} \rfloor (\text{the greatest integer } \leq r_j) & r_j \neq 0 \\ 0 & r_j = 0, \end{cases}$$

so  $r_{j+1} := r_j - q_j$  satisfies  $0 \leq r_{j+1} < 1$ .

If  $r_j = 0$  for all  $j > m$ , we also write the continued fraction expansion as the finite sequence  $[q_1, \dots, q_m] = q_1 + \frac{1}{q_2 + \frac{1}{q_3 + \dots + \frac{1}{q_m}}}$ , in which case we call the continued fraction expansion **finite**.

The  $q_j$  and  $r_j$  is used to make you think that these quantities are like quotients and remainders (which they are if  $x$  is rational). The rounding down function  $x \mapsto \lfloor x \rfloor$  (also often denoted by  $x \mapsto [x]$ ) is called the **greatest integer function** or the **floor function**.

Note for any  $x \in \mathbb{R}$ , there is a unique continued fraction expansion  $[q_1, q_2, \dots]$ . Moreover, since at each step  $0 \leq r_j < 1$ , the reciprocal will be at least 1 if  $r_j \neq 0$ , and so  $q_{j+1} = 0$  if and only if  $r_j = 0$ .

**Exercise 5.4.1.** Compute the continued fraction expansion of  $\frac{80}{17}$ .

**Exercise 5.4.2.** Let  $x \in \mathbb{R}$ , and  $[q_1, q_2, \dots]$  be the continued fraction expansion. Let  $(x_n)$  denote the sequence of rational numbers by evaluation the partial continued fraction expansions:

$$x_n = q_1 + \frac{1}{q_2 + \frac{1}{\dots + \frac{1}{q_n}}}$$

Show  $\lim_{n \rightarrow \infty} x_n = x$ .

**Exercise 5.4.3.** For  $x \in \mathbb{R}$ , show the continued fraction expansion for  $x$  is finite if and only if  $x \in \mathbb{Q}$ .

**Example 5.4.3.** Let's compute the continued fraction expansion  $[q_1, q_2, \dots]$  of  $\sqrt{5}$ .

First set

$$q_1 = \lfloor \sqrt{5} \rfloor = 2, \quad r_1 = \sqrt{5} - 2,$$

so at the first stage our expansion looks like

$$\sqrt{5} = q_1 + r_1 = 2 + (\sqrt{5} - 2) = 2 + \frac{1}{1/(\sqrt{5} - 2)}.$$

The nice thing about quadratic numbers is we can rationalize the denominator in  $\frac{1}{\sqrt{5}-2}$  by multiplying by the conjugate (in  $\mathbb{Z}[\sqrt{5}]$ ) of the denominator. Note  $N(r_1) = r_1 \bar{r}_1 = N(-2 + \sqrt{5}) = 4 - 5 = -1$ , so  $\frac{1}{r_1} = -\bar{r}_1 = 2 + \sqrt{5}$ , i.e.,

$$\frac{1}{r_1} = \frac{1}{\sqrt{5} - 2} = 2 + \sqrt{5}.$$

So at the next stage we let

$$q_2 = \lfloor 2 + \sqrt{5} \rfloor = 4, \quad r_2 = 2 + \sqrt{5} - q_2 = \sqrt{5} - 2 = r_1.$$

Thus at the next stage, we have the expansion

$$\sqrt{5} = 2 + \frac{1}{4 + \frac{1}{\sqrt{5}-2}}.$$

Since  $r_2 = r_1$ , we see that  $q_3 = q_2$ , so  $r_3 = r_2 = r_1$ , and so on. So these computations simply repeat, and we will have  $q_j = 4$  for all  $j \geq 2$ , giving the continued fraction expansion

$$\sqrt{5} = [2, 4, 4, 4, \dots] = 2 + \frac{1}{4 + \frac{1}{4 + \frac{1}{4 + \dots}}}.$$

**Example 5.4.4.** Now let's try finding the continued fraction expansion of  $\sqrt{3}$ . At the first stage we have

$$q_1 = \lfloor \sqrt{3} \rfloor = 1, \quad r_1 = \sqrt{3} - 1.$$

Then

$$\frac{1}{r_1} = \frac{1}{\sqrt{3} - 1} \frac{\sqrt{3} + 1}{\sqrt{3} + 1} = \frac{1 + \sqrt{3}}{2}.$$

So

$$q_2 = \lfloor \frac{1 + \sqrt{3}}{2} \rfloor = 1, \quad r_2 = \frac{1 + \sqrt{3}}{2} - 1 = \frac{\sqrt{3} - 1}{2}.$$

Then

$$\frac{1}{r_2} = \frac{2}{\sqrt{3} - 1} \frac{\sqrt{3} + 1}{\sqrt{3} + 1} = 1 + \sqrt{3}.$$

So

$$q_3 = \lfloor 1 + \sqrt{3} \rfloor = 2, \quad r_3 = \sqrt{3} - 1 = r_1.$$

Since  $r_3 = r_1$ , we must repeat after this, i.e.,  $q_4 = q_2$ ,  $r_4 = r_2$ , and so on, giving us the continued fraction expansion

$$\sqrt{3} = [1, 1, 2, 1, 2, 1, 2, 1, 2, \dots] = 1 + \frac{1}{1 + \frac{1}{2 + \frac{1}{1 + \frac{1}{2 + \dots}}}}$$

Notice the above continued fraction expansions repeat. We make this notion precise as follows.

**Definition 5.4.2.** We say a continued fraction expansion  $[q_1, q_2, \dots]$  is **periodic** if there exist  $s \geq 0$  and  $m \in \mathbb{N}$  such that

$$[q_1, q_2, \dots] = [q_1, \dots, q_s, q_{s+1}, \dots, q_{s+m}, q_{s+1}, \dots, q_{s+m}, q_{s+1}, \dots, q_{s+m}, \dots],$$

i.e., if  $q_{j+m} = q_j$  for all  $j > s$ . In this case, we denote this expansion by

$$[q_1, \dots, q_s, \overline{q_{s+1}, \dots, q_{s+m}}].$$

If the expansion is periodic, the smallest such  $m$  for which the above condition holds (for some  $s$ ) is called the **period** of  $[q_1, q_2, \dots]$ .

For instance, the examples above say  $\sqrt{5} = [2, \overline{4}]$  is periodic with period 1 and  $\sqrt{3} = [1, \overline{1, 2}]$  is periodic with period 2. Since any rational number has continued fraction expansion of the form  $[q_1, \dots, q_s, \overline{0}]$ , any rational number has a periodic continued fraction expansion with period 1. Note periodic continued fractions can be specified by a finite amount of data.

**Theorem 5.4.3** (Lagrange). For any  $x \in \mathbb{R}$ , the continued fraction expansion of  $x$  is periodic if and only if  $x \in \mathbb{Q}(\sqrt{d})$  for some  $d \geq 1$ . In particular the continued fraction expansion of any element of  $\mathbb{Z}[\sqrt{d}]$  is periodic.

Due to lack of time, we won't prove this. But the idea of the proof for the "if" direction is that at each stage in the continued fraction expansion, the quantities  $\frac{1}{r_j}$  will be elements of  $\mathbb{Q}(\sqrt{d})$  satisfying certain conditions, and then showing that there are only finitely many possibilities, so for some  $j, m \geq 1$ , we have  $r_{j+m} = r_j$  by the pigeonhole principle.

The "only if" direction is easier. For simplicity, we just illustrate the special case  $s = 0$ , so  $x = [\overline{q_1, \dots, q_m}]$  (what is called *purely periodic*). Consider  $\alpha = [q_1, \dots, q_m] \in \mathbb{Q}$ . Then

$$x = \alpha + \frac{1}{\alpha + \frac{1}{\alpha + \dots}}$$

Then

$$x - \alpha = \frac{1}{\alpha + \frac{1}{\alpha + \dots}},$$

so taking reciprocals shows

$$\frac{1}{x - \alpha} = \alpha + \frac{1}{\alpha + \frac{1}{\alpha + \dots}} = x,$$

i.e.,

$$1 = (x - \alpha)x = x^2 - \alpha x,$$

so  $x^2$  satisfies the quadratic equation  $x^2 - \alpha x - 1 = 0$  with rational coefficients, from which it follows  $x \in \mathbb{Q}(\sqrt{d})$  where  $d = \alpha^2 + 4$ .

**Exercise 5.4.4.** Compute the continued fraction expansion of  $\sqrt{2}$ .

**Exercise 5.4.5.** Compute the continued fraction expansion of  $\sqrt{7}$ .

Now we explain (without proof) the connection with Pell's equation and fundamental units. Assume  $d > 1$  is squarefree, and write

$$d = [q_1, \dots, q_s, \overline{q_{s+1}, \dots, q_{s+m}}]$$

where  $s$  and  $m$  are chosen minimally so we can represent this continued fraction periodically. In particular  $m$  is the period. Consider the partial continued fraction expansions  $[q_1, \dots, q_n]$ . These are rational numbers, so we write them as

$$\frac{x_n}{y_n} = [q_1, \dots, q_n], \quad x_n, y_n \in \mathbb{N}, \gcd(x_n, y_n) = 1.$$

As they converge to  $\sqrt{d}$  ([Exercise 5.4.2](#)), we call  $(x_n, y_n)$  the  $n$ -th **convergent** of the continued fraction.

**Theorem 5.4.4.** For  $d > 1$  squarefree, let  $(x_n, y_n)$  be the  $n$ -th convergent in the continued fraction expansion of  $\sqrt{d}$ . Then  $x_m + y_m\sqrt{d}$  is the fundamental unit  $\varepsilon_d$ , where  $m$  is the period of this continued fraction. More generally,  $x_{km} + y_{km}\sqrt{d}$  is  $\varepsilon_d^k$  for  $k \geq 1$ .

So in summary, we used units in real quadratic fields to determine all solutions to Pell's equation in terms of  $\varepsilon_d$ . Now we know how to compute  $\varepsilon_d$  in terms of continued fractions, and thus determine all solutions to Pell's equation. Moreover, this allows us to construct good rational approximations to  $\sqrt{d}$ . Of course, using continued fractions directly gives us rational approximations to  $\sqrt{d}$ , but in some sense the ones coming from solutions to Pell's equation (or  $x^2 - dy^2 = -1$ ) are optimal in that they will have minimal remainder (see [Proposition 5.2.3](#)). (We also haven't proved that the continued fraction convergents give us good rational approximations, though one can prove this.)

**Example 5.4.5.** Recall  $\sqrt{5} = [2, \overline{4}]$ , which has period 1. So we look at the first convergent is given by

$$\frac{x_1}{y_1} = [q_1] = [2] = \frac{2}{1}$$

so the theorem says  $\varepsilon_5 = 2 + \sqrt{5}$ , which matches with [Example 5.3.2](#).

**Example 5.4.6.** Recall  $\sqrt{2} = [1, \overline{1, 2}]$ , which has period 2. Thus the second convergent is given by

$$\frac{x_2}{y_2} = [1, 1] = 1 + \frac{1}{1} = \frac{2}{1},$$

so the theorem says  $\varepsilon_2 = 2 + \sqrt{2}$ , which matches what you should have got in [Exercise 5.3.1](#).

**Exercise 5.4.6.** Use the continued fraction expansion of  $\sqrt{7}$  to compute  $\varepsilon_7$  (see also [Example 5.3.3](#)). Then obtain a rational approximation to  $\sqrt{7}$  accurate to within  $\frac{1}{100}$ .

**Exercise 5.4.7.** Use the continued fraction expansion of  $\sqrt{19}$  to compute  $\varepsilon_{19}$ . Use  $\varepsilon_{19}$  to obtain a rational approximation to  $\sqrt{19}$  accurate to within  $\frac{1}{1000}$ . (You may use a calculator.)

**Exercise 5.4.8.** Use continued fractions to obtain the expression for  $\varepsilon_{61}^+$  asserted in [Example 5.4.1](#). (You may use a calculator.)

## 5.5 Aftermission: fundamental units and Fibonacci numbers

We close this chapter with an amusing connection with fundamental units and Fibonacci numbers, following ideas that we used to solve Pell's equation.

The **golden ratio**  $\phi = \frac{1+\sqrt{5}}{2}$  is the fundamental unit for the full ring of integers  $\mathbb{Z}[\frac{1+\sqrt{5}}{2}]$ . For  $x, y \in \mathbb{Z}$ , note

$$N(x + y\frac{1+\sqrt{5}}{2}) = (x + y\frac{1+\sqrt{5}}{2})(x + y\frac{1-\sqrt{5}}{2}) = x^2 + xy - y^2.$$

This expression is a binary quadratic form, which we also denote

$$Q(x, y) = x^2 + xy - y^2.$$

Recall the **Fibonacci numbers**  $F_n$  are defined by

$$F_1 = F_2 = 1, F_{n+2} = F_{n+1} + F_n, \quad n \geq 1.$$

**Exercise 5.5.1.** Show the Fibonacci numbers satisfy  $F_{2n+2}^2 + 1 = F_{2n+2}F_{2n+1} + F_{2n+1}^2$ .

Put another way, the exercise says that  $(F_{2n+1}, F_{2n+2})$  are solutions to

$$Q(x, y) = 1,$$

which is an analogous equation to Pell's equation. In other words

$$F_{2n+1} + F_{2n+2}\frac{1+\sqrt{5}}{2}$$

is a unit of norm 1 in  $\mathbb{Z}[\frac{1+\sqrt{5}}{2}]$ .

Note the golden ratio  $\phi$  has norm  $-1$ , but its square  $\varepsilon = \phi^2 = \frac{3+\sqrt{5}}{2}$  has norm 1 (i.e., is the fundamental +unit in  $\mathbb{Z}[\frac{1+\sqrt{5}}{2}]$ ). Then we can write

$$\varepsilon = \frac{3 + \sqrt{5}}{2} = 1 + 1 \cdot \frac{1 + \sqrt{5}}{2} = F_1 + F_2 \frac{1 + \sqrt{5}}{2}.$$

Computing a couple of powers of  $\varepsilon$ , we see

$$\varepsilon^2 = \frac{7 + 3\sqrt{5}}{2} = 2 + 3 \cdot \frac{1 + \sqrt{5}}{2} = F_3 + F_4 \frac{1 + \sqrt{5}}{2},$$

$$\varepsilon^3 = \frac{47 + 21\sqrt{5}}{2} = 13 + 21 \cdot \frac{1 + \sqrt{5}}{2} = F_5 + F_6 \frac{1 + \sqrt{5}}{2}.$$

This is part of a general rule.

**Exercise 5.5.2.** Compute  $\varepsilon^4$  directly and then check that  $\varepsilon_4 = F_7 + F_8 \frac{1+\sqrt{5}}{2}$ .

**Exercise 5.5.3.** Prove that  $\varepsilon^n = F_{2n-1} + F_{2n}\phi$  for  $n \geq 1$ .

The above expression gives a way to compute Fibonacci numbers. While it's not exactly presented as a formula for  $F_n$ , you probably noticed in the calculations above you immediately see  $F_{2n}$  as the coefficient  $b$  in the expression  $\varepsilon^n = \frac{a+b\sqrt{5}}{2}$ , and then  $a$  is just  $b + 2F_{2n-1}$ . One can rewrite these calculations into a well-known formula for  $F_n$ :

**Exercise 5.5.4.** Prove that  $F_n = \frac{\phi^n - \bar{\phi}^n}{\phi - \bar{\phi}}$  for  $n \geq 1$ .