

154: Diophantine Equations

Nir Elber

Fall 2023

CONTENTS

How strange to actually have to see the path of your journey in order to make it.

—Neal Shusterman, [Shu16]

Contents	2
1 Linear Equations	4
1.1 Modular Arithmetic and Sage	4
1.1.1 Local Obstructions	4
1.1.2 The Law of Linear Reciprocity	5
1.1.3 Bézout's Theorem	7
1.1.4 The Extended Euclidean Algorithm	8
1.1.5 Problems	10
1.2 Finite Continued Fractions	11
1.2.1 Connection to Continued Fractions	11
1.2.2 Continued Fraction Convergents	13
1.2.3 More on the Magic Box Algorithm	16
1.2.4 Problems	18
1.3 Infinite Continued Fractions	19
1.3.1 Convergence of Infinite Continued Fractions	19
1.3.2 Building Infinite Continued Fractions	22
1.3.3 Quadratic Irrationals	24
1.3.4 Convergents Are Good Rational Approximations	27
1.3.5 Convergents Are Best Rational Approximations	30
1.3.6 Problems	31
1.4 Diophantine Approximation	31
1.4.1 Irrationality Measure	32
1.4.2 Irrationality Measure via Continued Fractions	35
1.4.3 Algebraic Bounds on Irrationality Measure	37
1.4.4 e Is Transcendental	40
1.4.5 The Continued Fraction of e	44
1.4.6 Problems	47

2 Quadratic Equations	49
2.1 Pell Equations	49
2.1.1 Pell Equations via Elementary Methods	49
2.1.2 Pell Equations with Sophistication	53
2.1.3 Using Continued Fractions	55
2.1.4 Generalized Pell Equations	58
2.1.5 A Harder Problem	61
2.1.6 Problems	61
2.2 Number Rings	62
2.2.1 Normal Domains	63
2.2.2 Number Rings	64
2.2.3 The Discriminant	69
2.2.4 Number Ring Structure	69
2.2.5 Dirichlet's Unit Theorem: Upper Bound	69
2.3 Minkowski Theory	69
2.3.1 Minkowski's Theorem	69
2.3.2 Dirichlet's Unit Theorem: Upper Bound	69
2.4 Binary Quadratic Forms	69
3 Intermission: Other Fields	70
3.1 Cyclotomic Extensions	70
3.2 (Almost) Unique Factorization	70
3.3 Local Fields	70
3.4 Hensel's Lemma	70
4 Cubic Equations	71
4.1 Elliptic Curves	71
4.2 Torsion of Elliptic Curves	71
4.3 Elliptic Curves over Finite Fields	71
4.4 Modern Perspectives	71
A Some Algebra	72
A.1 Unique Factorization Domains	72
A.2 A Little Field Theory	76
A.2.1 Basic Notions	76
A.2.2 Polynomial Rings	77
A.2.3 Algebraic Elements	78
A.2.4 Enough Galois Theory to be Dangerous	81
A.2.5 Norm and Trace	82
Bibliography	83
List of Definitions	84

THEME 1

LINEAR EQUATIONS

Think deeply of simple things

—Ross Program, [Pro22]

1.1 Modular Arithmetic and Sage

In this section, we review the elementary number theory we will use in these notes. The goal of the present chapter is to be able to solve the equation

$$ax + by = 1$$

as quickly as possible, but we will encounter Diophantine approximation in the process.

1.1.1 Local Obstructions

A theme that will reappear in this course is that of “local obstructions,” so we introduce the idea now. Here are some examples.

Example 1.1. The only integer solution to the equation $x^2 + y^2 = 3z^2$ is $(x, y, z) = (0, 0, 0)$.

Solution. Of course $(0, 0, 0)$ is a solution, so the main content is showing that it is the only one. Suppose that (x, y, z) is a nonzero solution, and we suppose that (x, y, z) is minimal with respect to $|x| + |y| + |z| > 0$. If all the terms are even, then $(x/2, y/2, z/2)$ is also an integer solution with $|x/2| + |y/2| + |z/2| < |x| + |y| + |z|$, violating minimality. Thus, we may assume that at least one of the terms is odd. We have two cases; the main point is that $x^2 \equiv 0, 1 \pmod{4}$ for any integer x .

- If z is odd, then we are asking for

$$x^2 + y^2 \equiv 3 \pmod{4}.$$

But $x^2, y^2 \pmod{4} \in \{0, 1\}$ cannot achieve this.

- If z is even, then we are asking for

$$x^2 + y^2 \equiv 0 \pmod{4}.$$

However, without loss of generality we will have x odd and so $x^2 \equiv 1 \pmod{4}$. But then $x^2 + y^2 \equiv 1 + y^2 \pmod{4}$ will never be $0 \pmod{4}$.

All cases have caused contradiction, so we have finished the proof. ■

Example 1.2. There are no integer solutions to the equation $6x + 9y = 2$.

Solution. Reducing $(\text{mod } 3)$ means that any integer solution to $6x + 9y = 2$ implies $0 \equiv 2 \pmod{3}$, which is a contradiction. ■

Now that we've seen some examples, let's make explicit what is going on.



Idea 1.3. Given an equation $f(x_1, \dots, x_n) = 0$, we can check if f has solutions in \mathbb{Z} by first checking if there are solutions to

$$f(x_1, \dots, x_n) \equiv 0 \pmod{m}$$

for integers m .

What is useful about Idea 1.3 is that checking for solutions $(\text{mod } m)$ amounts to a finite computation where variables live in $\mathbb{Z}/m\mathbb{Z}$, and we can simply run the finite computation to check.

Of course, Idea 1.3 is not perfectly robust, but it will guide our discussion of Diophantine equations throughout this course.

Non-Example 1.4. One can show that

$$(x^2 - 2)(x^2 - 3)(x^2 - 6) = 0$$

has solutions $(\text{mod } p)$ for all primes p , but there is no integer solution.

Here is an example which is akin to Idea 1.3 but not quite the same.

Example 1.5. There are no integer solutions to $x^2 + y^2 = 2xy - 1$.

Solution. This equation is actually $(x - y)^2 = -1$, which has no solutions because $(x - y)^2 > -1$ for any real numbers $x, y \in \mathbb{R}$. ■

Example 1.6. There are no integer solutions to $x^2 + y^2 = 6$.

Solution. We see that $x \in \{0, \pm 1, \pm 2\}$ forces $y \in \{\pm\sqrt{6}, \pm\sqrt{5}, \pm\sqrt{2}\}$, none of which provide integer solutions. However, if $|x| \geq 3$, then

$$x^2 + y^2 = 9 + y^2 > 6,$$

from which we see that there are not even real solutions! ■

The above examples teach us that it is also useful to check for real-valued solutions to an equation in addition to checking $(\text{mod } m)$ for various integers m . These are also “local obstructions.”

1.1.2 The Law of Linear Reciprocity

Idea 1.3 is useful for determining when a linear equation of the form $ax + by = 1$ cannot have solutions. The goal of the present section is to show that these “local obstructions” are the only obstructions. Namely, we will prove a result of the following type.

Proposition 1.7. Let a, b , and c be integers. Then there are integers $x, y \in \mathbb{Z}$ such that $ax + by = c$ if and only if, for any integer m , there are integers $x_m, y_m \in \mathbb{Z}$ such that

$$ax_m + by_m \equiv c \pmod{m}.$$

In other words, it is enough to check locally. However, Proposition 1.7 is not very helpful for actually trying to determine if $ax + by = c$ has solutions: we would have to check $ax + by \equiv c \pmod{m}$ for infinitely many moduli m , which is not a finite computation! Thankfully, we have the following more effective version of Proposition 1.7.

Proposition 1.8. Let a, b , and c be integers. Then there are integers $x, y \in \mathbb{Z}$ such that $ax + by = c$ if and only if there are integers $x, y \in \mathbb{Z}$ such that

$$ax + by \equiv c \pmod{b}.$$

In other words, the only modulus we have to check is $m = b$. Let's prove Proposition 1.8.

Proof of Proposition 1.8. Of course having integers x and y such that $ax + by = c$ will imply that $ax + by \equiv c \pmod{b}$. Conversely, suppose we have integers x_0 and y_0 such that

$$ax_0 + by_0 \equiv c \pmod{b}.$$

Then we know there is some integer y_1 such that

$$ax_0 + by_0 = c + by_1,$$

so $ax_0 + b(y_0 - y_1) = c$ provides an integer solution to $ax + by = c$. ■

Example 1.9. The equation $3x + 5y = 1$ has integer solutions.

Solution. By Proposition 1.8, it suffices to check $\pmod{3}$. Then we are looking for integers x and y such that

$$3x + 5y \equiv 1 \pmod{3}.$$

Well, $(x, y) = (0, 2)$ will do the trick. ■

Example 1.10. The equation $2x + 4y = 3$ has no integer solutions.

Solution. By Proposition 1.8, it suffices to check $\pmod{2}$. Then we are looking for integers x and y such that

$$2x + 4y \equiv 3 \pmod{2}.$$

But this implies $0 \equiv 3 \pmod{2}$, which is a contradiction, so there can be no integer solutions. ■

Proposition 1.8 also allows us to prove the “reciprocity” theorem. These are also a major theme in number theory, though we will not see even close to the full story in this course. What is remarkable in the following result is that we have found a way to switch the modulus of our “local obstruction” around, perhaps at the cost of adjusting the equation being considered. Such statements are in general very profitable!

Proposition 1.11 (law of linear reciprocity). Let a, b , and c be integers. Then there is an integer x such that $ax \equiv c \pmod{b}$ if and only if there is an integer x such that $bx \equiv c \pmod{a}$.

Proof. There is an integer x such that $ax \equiv c \pmod{b}$ if and only if there are integers x and y such that $ax = c - by$, which is equivalent to

$$ax + by = c.$$

This condition is now symmetric in a and b , so running the above argument backwards provides equivalence to finding an integer x such that $bx \equiv c \pmod{a}$. ■

Example 1.12. The equation $93x + 35y = 1$ has integer solutions.

Solution. By Proposition 1.8, it is equivalent to check that

$$23x \equiv 93x + 35y \equiv 1 \pmod{35}$$

has integer solutions. By Proposition 1.11, this is equivalent to having integer solutions to

$$12x \equiv 35x \equiv 1 \pmod{23}.$$

Going again, by Proposition 1.11, this is equivalent to having integer solutions to

$$11x \equiv 23x \equiv 1 \pmod{12}.$$

Continuing, by Proposition 1.11, this is equivalent to having integer solutions to

$$x \equiv 12x \equiv 1 \pmod{11},$$

for which we see that $x = 1$ works. ■

Example 1.13. The equation $289x + 323y = 2$ has no integer solutions.

Solution. By Proposition 1.8, it is equivalent to check that

$$34y \equiv 289x + 323y \equiv 2 \pmod{289}$$

has integer solutions. By Proposition 1.11, this is equivalent to having integer solutions to

$$17x \equiv 289x \equiv 2 \pmod{34}.$$

One more time, Proposition 1.11 says that it is equivalent to have integer solutions to

$$0 \equiv 34x \equiv 2 \pmod{17},$$

which is false. ■

1.1.3 Bézout's Theorem

Proposition 1.11 does a good job of determining when there are integer solutions to an equation of the form $ax + by = c$, but we would like a more efficient characterization, and we would also like an efficient way to write down the solutions. We begin with the more uniform characterization.

Theorem 1.14 (Bézout). Let a , b , and c be integers. Then there are integers x and y such that $ax + by = c$ if and only if $\gcd(a, b)$ divides c .

We are going to prove Theorem 1.14 multiple times, essentially to emphasize different points of view on this area of number theory. To begin, let's establish that Proposition 1.11 is in fact able to provide a proof.

Proof of Theorem 1.14 via Proposition 1.11. We imitate the previous examples. Note that $ax + by = c$ if and only if $(-a)(-x) + by = c$ and similar for other choices of signs, so we might as well assume that a and b and c are all nonnegative integers. Additionally, having solutions for $ax + by = c$ is a condition symmetric on a and b , so we might as well assume that $a \leq b$.

We induct on a . If $a = 0$, then either $b = 0$, and we have a solution if and only if $c = 0 = \gcd(a, b)$, or $b \neq 0$, and we have a solution if and only if $c = by = \gcd(a, b)y$ for some integer y . Otherwise, $a > 0$. Now, by Proposition 1.8, we have an integer solution if and only if

$$ry \equiv ax + by \equiv c \pmod{a}$$

has an integer solution, where r is chosen so that $b \equiv r \pmod{a}$ and $0 \leq r < a$. By Proposition 1.11, this is now equivalent to having an integer solution to

$$ax \equiv c \pmod{b-a},$$

which by Proposition 1.8 is equivalent to having an integer solution to $rx + ay = c$. But now we have replaced (a, b) with (r, a) , where $r < a$ and $\gcd(a, b) = \gcd(r, a)$, so induction completes the argument. ■

The above argument is fairly involved, so it is rewarding to know that the following cleaner proof exists.

Proof of Theorem 1.14 via well-ordering. It suffices to show that

$$\{ax + by : x, y \in \mathbb{Z}\} = \gcd(a, b)\mathbb{Z}.$$

Quickly, if $a = b = 0$, then both sides are $\{0\}$, so there is nothing to say. Otherwise, we may assume that at least one of a or b is nonzero. Certainly $\gcd(a, b)$ divides $ax + by$ for any $x, y \in \mathbb{Z}$, so $\{ax + by : x, y \in \mathbb{Z}\} \subseteq \gcd(a, b)\mathbb{Z}$. It remains to show the other inclusion, which is equivalent to showing $\gcd(a, b) \in \{ax + by : x, y \in \mathbb{Z}\}$.

Well, we expect $\gcd(a, b)$ to be the smallest positive element of $\{ax + by : x, y \in \mathbb{Z}\}$, so we let g denote this smallest positive element, and we want to show that $g = \gcd(a, b)$. (This g exists by the well-ordering of \mathbb{N} . Note that $\{ax + by : x, y \in \mathbb{Z}\}$ certainly has some positive element because it contains $a^2 + b^2 > 0$.) Certainly $\gcd(a, b)$ divides g by the argument of the previous paragraph, so it suffices to show that g divides $\gcd(a, b)$, for which we will show that $g \mid a$ and $g \mid b$.

In fact, we will only show that $g \mid a$, and $g \mid b$ follows symmetrically. For this, we use the division algorithm to write

$$a = gq + r$$

for some integers $q, r \in \mathbb{Z}$ where $0 \leq r < g$. Now, $r = a - gq$ will live in $\{ax + by : x, y \in \mathbb{Z}\}$, but $r < g$ forces r to not be a positive element in this set by minimality, so we must have $r = 0$. Thus, $a = gq$, which means $g \mid a$, as needed. ■

The drawback of the above cleaner proof is that it is difficult to see how to turn it into an effective algorithm to actually compute x and y . Indeed, the argument does not even make it clear how to find $x, y \in \mathbb{Z}$ such that

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

which is in some sense the crux of the matter because we can then multiply x and y by $c/\gcd(a, b)$. With some care, we will be able to provide an effective algorithm, but it will take some care.

1.1.4 The Extended Euclidean Algorithm

The motivation to our algorithm will begin with wanting an efficient way to compute $\gcd(a, b)$, which we need to use Theorem 1.14 anyway. The Euclidean algorithm is based on the following lemma.

Lemma 1.15. Let a and b be integers. For any integer q , we have $\gcd(a, b) = \gcd(a - bq, b)$.

Proof. Note that an integer d divides a and b implies that d divides $a - bq$ and b ; the converse holds by a symmetric argument. Thus, the conclusion follows from taking the least elements of the sets

$$\{d \in \mathbb{Z}_{\geq 0} : d \mid a \text{ and } d \mid b\} = \{d \in \mathbb{Z}_{\geq 0} : d \mid a - bq \text{ and } d \mid b\},$$

finishing. ■

We are now equipped to see an example of the Euclidean algorithm.

Example 1.16. We use the “Euclidean algorithm” to compute $\gcd(93, 35)$.

Solution. To begin, we repeatedly use the division algorithm to write

$$\begin{aligned} 93 &= 2 \cdot 35 + 23 \\ 35 &= 1 \cdot 23 + 12 \\ 23 &= 1 \cdot 12 + 11 \\ 12 &= 1 \cdot 11 + 1 \\ 11 &= 11 \cdot 1 + 0. \end{aligned}$$

Thus, repeatedly applying Lemma 1.15, we see

$$\gcd(93, 35) = \gcd(35, 23) = \gcd(23, 12) = \gcd(12, 11) = \gcd(11, 1) = 1,$$

which is what we wanted. ■

Exercise 1.17. Use the Euclidean algorithm to compute $\gcd(47, 31)$.

It is somewhat technical to make a rigorous argument avoid the above process. Take a moment to read and digest the following statement.

Proposition 1.18 (Euclidean algorithm). Let a_0 and a_1 be positive coprime integers. Define the integer sequences a_2, a_3, \dots and q_0, q_1, \dots recursively by

$$a_n = q_n a_{n+1} + a_{n+2} \quad \text{where} \quad 0 \leq a_{n+2} < a_{n+1}$$

where $q_n := \lfloor a_n / a_{n+1} \rfloor$ if $a_{n+1} > 0$ and $(a_{n+2}, q_n) := (0, 0)$ otherwise. Then there is a minimal N such that $a_n = 0$ for $n > N$, and $a_N = \gcd(a_0, a_1)$.

Proof. By construction of the sequence, if $a_{n+1} > 0$, then $0 \leq a_{n+2} < a_{n+1}$. Thus, if $a_{n+1} > 0$ always, then a_1, a_2, \dots is a strictly decreasing sequence of positive integers, which is impossible by the well-ordering of the positive integers.

So indeed, there is some integer N such that $a_{N+1} = 0$, and we may choose N to be minimal with this property so that $a_N \neq 0$. (Note that $a_0 \neq 0$, so there is some n with $a_n \neq 0$.) Then $a_{N+1} = 0$ by construction, and the definition of our recursion enforces $a_n = 0$ for all $n > N$.

It remains to show that $a_N = \gcd(a_0, a_1)$. The main claim is that $\gcd(a_0, a_1) = \gcd(a_n, a_{n+1})$ for any $0 \leq n \leq N$, which will complete the proof by plugging in $n = N$. We show this claim by induction: there is nothing to say for $n = 0$, and for any $n < N$ so that $a_{n+1} > 0$, we see that

$$\gcd(a_n, a_{n+1}) = \gcd(q_n a_{n+1} + a_{n+2}, a_{n+1}) = \gcd(a_{n+1}, a_{n+2}),$$

which completes the inductive step. ■

Proposition 1.18 grants us another proof of Theorem 1.14.

Proof of Theorem 1.14 via Proposition 1.18. As usual, we start off with the “easier” direction: if $ax + by = c$ for some $x, y \in \mathbb{Z}$, then we note $\gcd(a, b)$ divides $ax + by$ and so divides c .

We use Proposition 1.18 to show the harder direction. Both the condition $ax + by = c$ and $\gcd(a, b) \mid c$ remain invariant to adjusting the sign of a and b , so we may assume $a, b \geq 0$. Additionally, if $a = 0$, then both conditions are equivalent to $b \mid c$; a symmetric argument works for $b = 0$. Thus, we may assume that $a, b > 0$.

Now, set $a_0 := a$ and $a_1 := b$ and build the sequence a_2, a_3, \dots of Proposition 1.18. By induction, we see that

$$a_n \in \{a_0x + a_1y : x, y \in \mathbb{Z}\}.$$

Indeed, there is nothing to say for $n = 0$ and $n = 1$. Then for the induction, we note that $\{a_0x + a_1y : x, y \in \mathbb{Z}\}$ is closed under \mathbb{Z} -linear combination, so containing a_n and a_{n+1} implies containing $a_{n+2} = a_n - q_n a_{n+1}$. Thus, using Proposition 1.18, we see that $a_N = \gcd(a, b)$ takes the form $ax + by$ for $x, y \in \mathbb{Z}$, completing the proof. ■

We are finally able to read the above proof closely to have an effective algorithm to compute x and y solving $ax + by = \gcd(a, b)$. This is called the “extended Euclidean algorithm” and is best seen by example.

Example 1.19. We use the “extended Euclidean algorithm” to find integers x and y such that $93x + 35y = 1$.

Proof. The idea is to run the Euclidean algorithm backwards “solving” for the remainders. Indeed, using the computations of Example 1.16, we see

$$\begin{aligned} 1 &= 12 - 1 \cdot 11 \\ 11 &= 23 - 1 \cdot 12 \\ 12 &= 35 - 1 \cdot 23 \\ 23 &= 93 - 2 \cdot 35. \end{aligned}$$

We now plug in for each successive remainder, writing

$$\begin{aligned} 1 &= 12 - 1 \cdot 11 \\ &= 12 - 1 \cdot (23 - 1 \cdot 12) = 2 \cdot 12 - 1 \cdot 23 \\ &= 2 \cdot (35 - 1 \cdot 23) - 1 \cdot 23 = 2 \cdot 35 - 3 \cdot 23 \\ &= 2 \cdot 35 - 3 \cdot (93 - 2 \cdot 35) = 8 \cdot 35 - 3 \cdot 93. \end{aligned}$$

Thus, $(x, y) = (-3, 8)$ will do the trick. ■

Exercise 1.20. Use the extended Euclidean algorithm to find integers x and y such that $47x + 31y = 1$.

1.1.5 Problems

Do at least ten points worth of the following exercises.

Problem 1.1.1 (1 point). Let $n \equiv 3 \pmod{4}$. Show that there are not two integers $x, y \in \mathbb{Z}$ such that $x^2 + y^2 = n$.

Problem 1.1.2 (2 points). Let $n \equiv 7 \pmod{8}$. Show that there are not three integers $x, y, z \in \mathbb{Z}$ such that $x^2 + y^2 + z^2 = n$.

Problem 1.1.3 (2 points). Let a and b be integers. Suppose that there are pairs of integers (x, y) and (x', y') such that $ax + by = ax' + by' = 1$. Show that

$$x \equiv x' \pmod{b} \quad \text{and} \quad y \equiv y' \pmod{a}.$$

Problem 1.1.4 (2 points). Define the Fibonacci sequence $\{F_n\}_{n=0}^{\infty}$ by $F_0 = 0$, $F_1 = 1$, and $F_{n+2} = F_{n+1} + F_n$ for any $n \geq 0$. Show that $\gcd(F_{n+1}, F_n) = 1$ for any $n \geq 0$.

Problem 1.1.5 (3 points). Compute $\gcd(1027, 1738)$. Then find integers x and y such that $1027x + 1738y = \gcd(1027, 1738)$.

Problem 1.1.6 (3 points). Let a , b , and c be integers with $\gcd(a, b, c) = 1$. Show that there exist integers $x, y, z \in \mathbb{Z}$ such that $ax + by + cz = 1$.

Problem 1.1.7 (5 or 6 points). Implement the extended Euclidean algorithm.

(a) For five points, write (and submit) a function in Python which takes as input two coprime positive integers a and b and outputs integers x and y such that $ax + by = 1$. Your function should implement the extended Euclidean algorithm.

(b) For an additional point, make the function work for any coprime integers a and b .

Your test case is $(a, b) = (12345678901, 10987654321)$.

1.2 Finite Continued Fractions

In this section, we begin our discussion of continued fractions with a discussion of finite continued fractions. The reward for our efforts will be a more memory-efficient version of the extended Euclidean algorithm.

1.2.1 Connection to Continued Fractions

We begin with the definition of a continued fraction.

Definition 1.21 (continued fraction). A *continued fraction* expansion is an expression of the form

$$a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{a_3 + \ddots}}},$$

which we will notate by $[a_0; a_1, a_2, \dots]$. The terms a_i are the *continued fraction coefficients*.

In our application, the terms a_0, a_1, a_2, \dots will always be integers, and a_1, a_2, \dots will always be positive integers, but we take the moment to remark that this definition operates just fine even if these are not integers. This specialization does guarantee that we never run into division-by-zero problems, which is its principal advantage.

Remark 1.22. For the present section, our continued fractions will always be finite in length. In other words, our continued fractions will look like $[a_0; a_1, a_2, \dots, a_n]$ for some perhaps large n . In the next section, we will allow continued fractions to have infinite length by defining

$$[a_0; a_1, a_2, \dots] := \lim_{n \rightarrow \infty} [a_0; a_1, a_2, \dots, a_n],$$

but we will have to prove that this limit exists before providing this definition.

Continued fractions will be very interesting to us in the sequel, approximately speaking because they provide good rational approximations to real numbers. To start us off, suppose we have a real number α , and we would like to find coefficients $a_0, a_1, a_2, \dots \in \mathbb{Z}$ such that $\alpha = [a_0; a_1, a_2, \dots]$. In fact, we will be able to enforce $a_1, a_2, \dots \in \mathbb{Z}_{\geq 0}$. To see how, note that if we want

$$\alpha = a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{a_3 + \ddots}}},$$

then we should have $a_0 := \lfloor \alpha \rfloor$. Once we agree what a_0 should be, we may rearrange this equation into

$$\frac{1}{\alpha - a_0} = a_1 + \frac{1}{a_2 + \frac{1}{a_3 + \ddots}}.$$

Now we are trying to compute the continued fraction for $(\alpha - \lfloor \alpha \rfloor)^{-1} > 1$, so we may recurse. Namely, set $a_1 := \lfloor (\alpha - \lfloor \alpha \rfloor)^{-1} \rfloor$ and then rearrange again.

Here's an example.

Example 1.23. We express $93/35$ as a continued fraction.

Solution. We write

$$\begin{aligned} \frac{93}{35} &= 2 + \frac{23}{35} \\ &= 2 + \frac{1}{35/23} \\ &= 2 + \frac{1}{1 + \frac{12}{23}} \\ &= 2 + \frac{1}{1 + \frac{1}{23/12}} \\ &= 2 + \frac{1}{1 + \frac{1}{1 + \frac{11}{12}}} \\ &= 2 + \frac{1}{1 + \frac{1}{1 + \frac{1}{12/11}}} \\ &= 2 + \frac{1}{1 + \frac{1}{1 + \frac{1}{1 + \frac{1}{11}}}}, \end{aligned}$$

so $\frac{93}{35} = [2; 1, 1, 1, 11]$. ■

Exercise 1.24. Express $47/31$ as a continued fraction.

Compare Example 1.16 with Example 1.23: the coefficients $[2; 1, 1, 1, 11]$ match up exactly with the quotients appearing in the Euclidean algorithm. Rigorizing this is a little technical, but it is not too hard.

Proposition 1.25. Let a_0 and a_1 be coprime positive integers, and define integer sequences q_0, q_1, \dots, q_N and $a_0, a_1, a_2, \dots, a_N$ recursively as in Proposition 1.18 by

$$a_n = q_n a_{n+1} + a_{n+2}$$

for any $n \geq 0$, where $0 < a_{n+2} < a_{n+1}$ and terminating once $a_N = 1$ so that $a_{N+1} = 0$. Then $\frac{a_0}{a_1} = [q_0; q_1, q_2, \dots, q_N]$.

Proof. Recall that N exists by the Euclidean algorithm. We induct on N . If $N = 1$, then $a_1 = 1$ and

$$a_0 = q_0 a_1 + a_2$$

forces $a_2 = 0$ and $q_0 = a_0$. Thus, $a_0 = \frac{a_0}{a_1} = q_0 = [q_0]$.

Now take $N > 1$ (which implies $a_2 > 0$), and suppose the statement is true at $N - 1$. Then we see $a_0 = q_0 a_1 + a_2$ implies

$$\frac{a_0}{a_1} = q_0 + \frac{1}{a_1/a_2}.$$

Thus, running the Euclidean algorithm with the coprime positive integers a_1 and a_2 , we find that $\frac{a_1}{a_2} = [q_1; q_2, \dots, q_N]$ by the inductive hypothesis. It follows

$$\frac{a_0}{a_1} = q_0 + \frac{1}{[q_1; q_2, \dots, q_N]} = [q_0; q_1, q_2, \dots, q_N],$$

which is what we wanted. ■

Remark 1.26. Proposition 1.25 also has the nice side effect of showing that any rational number is equal to some finite continued fraction. However, note this continued fraction is not unique: given integers $a_0, a_1, a_2, \dots, a_n$ with a_1, a_2, \dots, a_n positive, one has

$$[a_0; a_1, a_2, \dots, a_{n-1}, a_n] = [a_0; a_1, a_2, \dots, a_{n-1}, a_n - 1, 1]$$

when $a_n > 1$, and otherwise

$$[a_0; a_1, a_2, \dots, a_{n-1}, 1] = [a_0; a_1, a_2, \dots, a_{n-1} + 1].$$

In particular, given any rational number q , we can find n of any parity such that there are integers $a_0, a_1, a_2, \dots, a_n$ with a_1, a_2, \dots, a_n positive and $q = [a_0; a_1, a_2, \dots, a_n]$.

The proof of Proposition 1.25 is fairly instructive: many of our arguments involving continued fractions are going to be inductive ones using identities like

$$q_0 + \frac{1}{[q_1; q_2, \dots, q_N]} = [q_0; q_1, q_2, \dots, q_N].$$

1.2.2 Continued Fraction Convergents

We mentioned at the outset that continued fractions provide good rational approximations for numbers. The way that this is done is by taking a long continued fraction $[a_0; a_1, a_2, \dots]$ and “truncating” it at some point to produce the shorter (and notably finite) continued fraction $[a_0; a_1, a_2, \dots, a_n]$. This truncation process is so important it has a name.

Definition 1.27 (convergent). Given a continued fraction $[a_0; a_1, a_2, \dots]$ and some $n \geq 0$, the truncation $[a_0; a_1, a_2, \dots, a_n]$ is the n th convergent, often denoted

$$\frac{h_n}{k_n} := [a_0; a_1, \dots, a_n].$$

As usual, we begin with an example.

Example 1.28. We compute the continued fraction convergents of $93/35$.

Solution. In Example 1.23, we computed that $\frac{93}{35} = [2; 1, 1, 1, 11]$, so here are our convergents.

- The zeroth convergent is $[2] = 2$.
- The first convergent is $[2; 1] = 2 + \frac{1}{1} = 3$.
- The second convergent is $[2; 1, 1] = 2 + \frac{1}{1+1} = \frac{5}{2}$.
- The third convergent is $[2; 1, 1, 1]$ is

$$[2; 1, 1, 1] = 2 + \frac{1}{1 + \frac{1}{1+1}} = 2 + \frac{1}{3/2} = \frac{8}{3}.$$

- The fourth convergent is $[2; 1, 1, 1, 11] = \frac{93}{35}$. ■

Exercise 1.29. Compute the continued fraction convergents of $47/31$.

The process outlined in Example 1.28 is rather annoying to execute by hand. We did not even compute $[2; 1, 1, 1, 11]$ by hand, but already $[2; 1, 1, 1]$ required some focus. In general, the problem with computing these convergents is that we are basically doing a totally new computation for every convergent.

However, there is a much faster way to compute these convergents: the “magic box” algorithm. For a sense of wonder, we will describe the algorithm first and then prove that it works second. We begin with the following grid.

$$\begin{array}{cc|ccccc} & & 2 & 1 & 1 & 1 & 11 \\ \hline 0 & 1 & & & & & \\ 1 & 0 & & & & & \end{array}$$

Explicitly, the 0s and 1s on the leftmost two columns will always be there in all computations, and the top row is made of our coefficients $[2; 1, 1, 1, 11]$. We now fill in the grid column-by-column, moving from left to right. For the next leftmost column, we multiply the coefficient 2 by the previous column and then add the column before that. In other words, we compute

$$\begin{bmatrix} 0 \\ 1 \end{bmatrix} + 2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix},$$

so the next column in our grid is as follows.

$$\begin{array}{cc|ccccc} & & 2 & 1 & 1 & 1 & 11 \\ \hline 0 & 1 & 2 & & & & \\ 1 & 0 & 1 & & & & \end{array}$$

Indeed, $2/1$ is the zeroth convergent. We now repeat the process: multiply 1 by the previous column and then add the column before that, writing

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix} + 1 \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \end{bmatrix},$$

making our grid as follows.

$$\begin{array}{cc|ccccc} & & 2 & 1 & 1 & 1 & 11 \\ 0 & 1 & 2 & 3 & & & \\ 1 & 0 & 1 & 2 & & & \end{array}$$

Indeed, $3/1$ is the first convergent. We now fill in the rest of the grid.

$$\begin{array}{cc|ccccc} & & 2 & 1 & 1 & 1 & 11 \\ 0 & 1 & 2 & 3 & 5 & 8 & 93 \\ 1 & 0 & 1 & 1 & 2 & 3 & 35 \end{array}$$

And indeed, we see the remaining convergents $5/2$, $8/3$, and $93/35$ appear from our grid.

Exercise 1.30. Execute this “magic box” algorithm to compute the continued fraction convergents of $47/31$.

Exercise 1.31. Compute the following 2×2 “minors” of our grid, as follows.

$$\det \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \det \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix}, \quad \det \begin{bmatrix} 2 & 3 \\ 1 & 1 \end{bmatrix}, \quad \det \begin{bmatrix} 3 & 5 \\ 1 & 2 \end{bmatrix}, \quad \dots$$

Do you see any patterns?

Proving that the magic box algorithm works is again somewhat technical. Perhaps the main difficulty is figuring out how to state the result, but the proof is still tricky. For now, we will settle for the following statement, but we will establish the refinement Corollary 1.36 shortly.

Proposition 1.32 (magic box). Let a_0, a_1, a_2, \dots be real numbers, where a_1, a_2, \dots are positive. Define the sequences $\{h_n\}_{n=-2}^\infty$ and $\{k_n\}_{n=-2}^\infty$ of real numbers recursively by

$$\begin{bmatrix} h_{-2} & h_{-1} \\ k_{-2} & k_{-1} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} h_{n+2} \\ k_{n+2} \end{bmatrix} = a_{n+2} \begin{bmatrix} h_{n+1} \\ k_{n+1} \end{bmatrix} + \begin{bmatrix} h_n \\ k_n \end{bmatrix}$$

for $n \geq -2$. Then

$$[a_0; a_1, \dots, a_n] = \frac{h_n}{k_n}$$

for any $n \geq 0$.

Proof. The requirement that a_1, a_2, \dots be positive is entirely to avoid division by zero errors. We also take a moment to recognize that the a_\bullet are being allowed to be real numbers rather than only integers. This will actually be relevant to the proof!

We induct on n . For $n = 0$, we can compute that $(h_0, k_0) = a_0(1, 0) + (0, 1) = (a_0, 1)$, so $\frac{h_0}{k_0} = a_0 = [a_0]$. For $n = 1$, we can compute that $(h_1, k_1) = a_1(a_0, 1) + (1, 0) = (a_1 a_0 + 1, a_1)$, so

$$\frac{h_1}{k_1} = \frac{a_1 a_0 + 1}{a_1} = a_0 + \frac{1}{a_1} = [a_0; a_1].$$

Now take $n \geq 2$. The trick for the inductive step is to write

$$[a_0; a_1, \dots, a_{n-2}, a_{n-1}, a_n] = a_0 + \frac{1}{a_1 + \frac{1}{\ddots + a_{n-2} + \frac{1}{a_{n-1} + \frac{1}{a_n}}}} = \left[a_0; a_1, \dots, a_{n-2}, a_{n-1} + \frac{1}{a_n} \right].$$

We may now apply the inductive hypothesis to this altered continued fraction, which is legal because $a_{n-1} + 1/a_n$ is surely a positive real number. Explicitly, define the sequence $a'_0, a'_1, \dots, a'_{n-1}$ where $a'_m := a_m$ for $m < n-1$ and $a'_{n-1} := a_{n-1} + \frac{1}{a_n}$, and then define the sequence $\{h'_m\}_{m=-2}^{n-1}$ and $\{k'_m\}_{m=-2}^{\infty}$ as in the proposition so that

$$[a_0; a_1, \dots, a_{n-2}, a_{n-1}, a_n] = [a'_0; a'_1, \dots, a'_{n-1}] = \frac{h'_{n-1}}{k'_{n-1}}.$$

To compute h'_{n-1} and k'_{n-1} we acknowledge that the construction of the a'_\bullet implies that $h'_m = h_m$ and $k'_m = k_m$ for $m < n-1$. So we see that

$$\begin{aligned} \begin{bmatrix} h'_{n-1} \\ k'_{n-1} \end{bmatrix} &= a'_{n-1} \begin{bmatrix} h'_{n-2} \\ k'_{n-2} \end{bmatrix} + \begin{bmatrix} h'_{n-3} \\ k'_{n-3} \end{bmatrix} \\ &= \left(a_{n-1} + \frac{1}{a_n} \right) \begin{bmatrix} h_{n-2} \\ k_{n-2} \end{bmatrix} + \begin{bmatrix} h_{n-3} \\ k_{n-3} \end{bmatrix} \\ &= \begin{bmatrix} \left(a_{n-1} + \frac{1}{a_n} \right) h_{n-2} + h_{n-3} \\ \left(a_{n-1} + \frac{1}{a_n} \right) k_{n-2} + k_{n-3} \end{bmatrix}. \end{aligned}$$

From here, we compute

$$\begin{aligned} \frac{h'_{n-1}}{k'_{n-1}} &= \frac{a_{n-1}a_n h_{n-2} + h_{n-2} + a_n h_{n-3}}{a_{n-1}a_n k_{n-2} + k_{n-2} + a_n k_{n-3}} \\ &= \frac{a_n(a_{n-1}h_{n-2} + h_{n-3}) + h_{n-2}}{a_n(a_{n-1}k_{n-2} + k_{n-3}) + k_{n-2}} \\ &= \frac{a_n h_{n-1} + h_{n-2}}{a_n k_{n-1} + k_{n-2}} \\ &= \frac{h_n}{k_n}, \end{aligned}$$

which completes the proof. ■

Remark 1.33. The proof of Proposition 1.32 in fact works even if we merely assume that the a_\bullet are indeterminate variables.

Example 1.34. Define the Fibonacci sequence $\{F_n\}_{n=0}^{\infty}$ by $F_0 = 0$ and $F_1 = 1$ and $F_{n+2} = F_{n+1} + F_n$ for any $n \geq 0$. Then for any $n \geq 0$,

$$\underbrace{[1; 1, \dots, 1]}_{n+1} = \frac{F_{n+2}}{F_{n+1}}.$$

Solution. We proceed by induction on n , using Proposition 1.32. From there, we may compute that $h_0/k_0 = 1/1 = F_2/F_1$ and $h_1/k_1 = 2/1 = F_3/F_2$. For the inductive step, we note that Proposition 1.32 yields

$$h_{n+2} = h_{n+1} + h_n \quad \text{and} \quad k_{n+2} = k_{n+1} + k_n$$

for any $n \geq 0$, which is the recursion for the Fibonacci sequence. ■

1.2.3 More on the Magic Box Algorithm

Proposition 1.32 essentially explains why the magic box works, though perhaps there is some doubt that the fractions h_n/k_n is in reduced form. Let's show this. We begin by explaining Exercise 1.31.

Corollary 1.35. Let a_0, a_1, a_2, \dots be real numbers, where a_1, a_2, \dots are positive, and define $\{h_n\}_{n=-2}^\infty$ and $\{k_n\}_{n=-2}^\infty$ as in Proposition 1.32. Then

$$\det \begin{bmatrix} h_n & h_{n+1} \\ k_n & k_{n+1} \end{bmatrix} = (-1)^{n+1}$$

for any $n \geq -2$.

Proof. This is essentially row-reduction. We proceed by induction on n . At $n = -2$, we see that $\det \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} = -1$. For the inductive step, suppose the statement for n , and we show $n + 1$. We note

$$\begin{bmatrix} h_{n+2} \\ k_{n+2} \end{bmatrix} = a_{n+2} \begin{bmatrix} h_{n+1} \\ k_{n+1} \end{bmatrix} + \begin{bmatrix} h_n \\ k_n \end{bmatrix}$$

allows us to use column operations in order to see

$$\det \begin{bmatrix} h_{n+1} & h_{n+2} \\ k_{n+1} & k_{n+2} \end{bmatrix} = \det \begin{bmatrix} h_{n+1} & h_n \\ k_{n+1} & k_n \end{bmatrix} = -\det \begin{bmatrix} h_n & h_{n+1} \\ k_n & k_{n+1} \end{bmatrix} = -(-1)^{n+1} = (-1)^{n+2},$$

which is what we wanted. ■

Corollary 1.36. Let a_0, a_1, a_2, \dots be integers, where a_1, a_2, \dots are positive, and define $\{h_n\}_{n=-2}^\infty$ and $\{k_n\}_{n=-2}^\infty$ as in Proposition 1.32. Then, for any $n \geq 0$,

$$[a_0; a_1, \dots, a_n] = \frac{h_n}{k_n},$$

and h_n/k_n is a fraction in reduced form with $k_n \geq 1$.

Proof. The equality follows directly from Proposition 1.32. Additionally, note that h_n and k_n are integers because they are terms of a sequence defined by integer recursion. Thus, to complete the proof, we must show that $\gcd(h_n, k_n) = 1$ and that $k_n \geq 1$ for $n \geq 0$. On one hand, we see $\gcd(h_n, k_n) = 1$ is direct from Corollary 1.35. On the other hand, $k_n \geq 1$ follows from a quick induction because $k_{-1} = 0$ and $k_0 = a_1 \geq 1$ and so $k_{n+2} = a_{n+2}k_{n+1} + k_n \geq 1$ always. ■

Corollary 1.35 has in fact suggested a faster algorithm (in terms of memory) than the Extended Euclidean algorithm. Let's see this by example.

Example 1.37. We find integers x and y such that $93x + 35y = 1$.

Solution. As in Example 1.16, we begin by writing

$$\begin{aligned} 93 &= 2 \cdot 35 + 23 \\ 35 &= 1 \cdot 23 + 12 \\ 23 &= 1 \cdot 12 + 11 \\ 12 &= 1 \cdot 11 + 1 \\ 11 &= 11 \cdot 1 + 0. \end{aligned}$$

From here, we apply the magic box algorithm Proposition 1.32 to build the following grid.

		2	1	1	1	11
0	1	2	3	5	8	93
1	0	1	1	2	3	35

Tracking Corollary 1.35 through, we see that

$$35 \cdot 8 - 93 \cdot 3 = \det \begin{bmatrix} 8 & 93 \\ 3 & 25 \end{bmatrix} = 1,$$

so $(x, y) = (-3, 8)$ works. ■

Remark 1.38. Here are a few ways to “check” the magic box algorithm.

- If using the magic box algorithm to compute convergents of the fraction p/q , then the last column of the magic box grid should yield p/q .
- The magic box algorithm has 2×2 minors controlled by Corollary 1.35, so one can compute a few of these for security.

1.2.4 Problems

Do at least 10 points worth of the following exercises.

Problem 1.2.1 (1 point). Find integer sequences $a_0, a_1, a_2, \dots, a_m$ and $b_0, b_1, b_2, \dots, b_n$ with a_1, a_2, \dots, a_m and b_1, b_2, \dots, b_n positive such that the sequences are distinct, but

$$[a_0; a_1, \dots, a_m] = [b_0; b_1, \dots, b_n].$$

Problem 1.2.2 (2 points). Compute the continued fraction convergents of $1738/1027$.

Problem 1.2.3 (3 points). Let a_0, a_1, a_2, \dots be integers, where a_1, a_2, \dots are positive, and define $\{h_n\}_{n=-2}^{\infty}$ and $\{k_n\}_{n=-2}^{\infty}$ as in Proposition 1.32. Show that

$$\left| \det \begin{bmatrix} h_n & h_{n+2} \\ k_n & k_{n+2} \end{bmatrix} \right| = |a_{n+2}|$$

for any $n \geq -2$. Additionally, predict the sign as a function on n .

Problem 1.2.4 (5 or 6 points). Let $a_0, a_1, a_2, \dots, a_m$ and $b_0, b_1, b_2, \dots, b_n$ be integers with a_1, a_2, \dots, a_m and b_1, b_2, \dots, b_n positive. Suppose

$$[a_0; a_1, a_2, \dots, a_m] = [b_0; b_1, b_2, \dots, b_n].$$

- For five points, suppose $m = n$. Show that $a_k = b_k$ for all $0 \leq k \leq m$.
- For an additional point, suppose $m < n$. Show that $m = n - 1$ and $a_k = b_k$ for $0 \leq k \leq m - 1$.

Problem 1.2.5 (5 points). Write (and submit) a function in Python which takes as input a list of integers $[a_0, a_1, a_2, \dots]$ with a_1, a_2, \dots positive and an index n and outputs the n th convergent $[a_0; a_1, a_2, \dots, a_n]$. You should implement the magic box algorithm.

Your test case is $[2; 1, 2, 1, 1, 4, 1, 1, 6, 1]$.

1.3 Infinite Continued Fractions

In this section, we examine continued fractions more closely. Our main task will be to show that continued fractions provide good and in fact the best rational approximations for a given irrational number. Of course, it will be a nontrivial task in order to make sense of what “best” means in this context. To set up our intuition, we will say that a fraction h/k provides a good rational approximation for a real number α if the difference

$$\left| \alpha - \frac{h}{k} \right|$$

is smaller than one might expect it to be. Of course, for any given denominator, we know that $[k\alpha] \leq k\alpha < [k\alpha] + 1$, so

$$\left| \alpha - \frac{[k\alpha]}{k} \right| \leq \frac{1}{k},$$

so a bound of $1/k$ is not too impressive. In fact, if α is irrational, we will be able to show that there are infinitely many rational numbers h/k such that

$$\left| \alpha - \frac{h}{k} \right| < \frac{1}{\sqrt{5}k^2},$$

and we will be able to show that this bound is essentially sharp.

1.3.1 Convergence of Infinite Continued Fractions

Thus far our discussion has been focused on finite continued fractions. We would now like to extend this discussion to infinite continued fractions. As in Remark 1.22, we would like to define

$$[a_0; a_1, a_2, \dots] \stackrel{?}{=} \lim_{n \rightarrow \infty} [a_0; a_1, a_2, \dots, a_n],$$

but we should begin by showing that this limit in fact exists. The idea is to show that the infinite continued fraction is an infinite series, and then we can use known results on infinite series to complete the proof. As such, we begin by turning $[a_0; a_1, a_2, \dots]$ into a series.

Lemma 1.39. Let a_0, a_1, a_2, \dots be real numbers, where a_1, a_2, \dots are positive, and let $\{h_n/k_n\}_{n=0}^\infty$ denote the continued fraction convergents $h_n/k_n := [a_0; a_1, a_2, \dots, a_n]$ where $k_n \geq 1$ and $\gcd(h_n, k_n) = 1$. Then

$$\frac{h_n}{k_n} - \frac{h_{n+1}}{k_{n+1}} = \frac{(-1)^{n+1}}{k_n k_{n+1}}.$$

Thus,

$$\frac{h_n}{k_n} = \frac{h_0}{k_0} + \sum_{m=0}^{n-1} \frac{(-1)^m}{k_m k_{m+1}}.$$

Proof. Note that $\{h_n\}_{n=0}^\infty$ and $\{k_n\}_{n=0}^\infty$ are the sequences constructed in Proposition 1.32 by Corollary 1.36. As such, the first claim follows directly from Corollary 1.35. The second claim now follows from writing

$$\frac{h_n}{k_n} = \frac{h_0}{k_0} + \sum_{m=0}^{n-1} \left(\frac{h_{m+1}}{k_{m+1}} - \frac{h_m}{k_m} \right) = \frac{h_0}{k_0} + \sum_{m=0}^{n-1} \frac{(-1)^m}{k_m k_{m+1}},$$

which is what we wanted. ■

Proposition 1.40. Let a_0, a_1, a_2, \dots be integers, where a_1, a_2, \dots are positive, and let $\{h_n/k_n\}_{n=0}^\infty$ denote the continued fraction convergents $h_n/k_n := [a_0; a_1, \dots, a_n]$ where $k_n \geq 1$ and $\gcd(h_n, k_n) = 1$. Then

$$\alpha := \lim_{n \rightarrow \infty} [a_0; a_1, a_2, \dots, a_n]$$

converges, and

$$\frac{1}{k_n(k_{n+1} + k_n)} < \left| \alpha - \frac{h_n}{k_n} \right| < \frac{1}{k_n k_{n+1}}$$

for each $n \geq 0$.

Proof. As usual, note that $\{h_n\}_{n=0}^\infty$ and $\{k_n\}_{n=0}^\infty$ are the sequences constructed in Proposition 1.32 by Corollary 1.36. To begin, we compute the limit as

$$\alpha = \lim_{n \rightarrow \infty} \frac{h_n}{k_n} = \frac{h_0}{k_0} + \sum_{n=0}^{\infty} \frac{(-1)^n}{k_n k_{n+1}},$$

where we have used Lemma 1.39 in the last equality. Now, the sequence $\{k_n\}_{n=0}^\infty$ is strictly increasing by Proposition 1.32 because a_1, a_2, \dots are all positive integers. Thus, the summation above absolute converges: an induction shows $k_n \geq n + 1$, so

$$\frac{h_0}{k_0} + \sum_{n=0}^{\infty} \left| \frac{(-1)^n}{k_n k_{n+1}} \right| \leq \frac{h_0}{k_0} + \sum_{n=0}^{\infty} \frac{1}{(n+1)(n+2)} < \infty.$$

As such, the limit does in fact converge.

To compute the error term, we use the error bound for alternating series. To begin the computation, note that the above work allows us to write

$$\left| \alpha - \frac{h_n}{k_n} \right| = \left| \frac{h_0}{k_0} + \sum_{m=0}^{\infty} \frac{(-1)^m}{k_m k_{m+1}} - \frac{h_0}{k_0} - \sum_{m=0}^{n-1} \frac{(-1)^m}{k_m k_{m+1}} \right| = \left| \sum_{m=n}^{\infty} \frac{(-1)^m}{k_m k_{m+1}} \right|.$$

Because the sequence $\{k_m\}_{m=0}^\infty$ is strictly increasing, the terms in the sum are monotonously decreasing in magnitude to zero, so the error bound for alternating series forces $|\alpha - h_n/k_n| < 1/(k_n k_{n+1})$, which proves the upper bound for our error.

To prove the lower bound of the error, we adjust for signs and note that the sum is

$$\begin{aligned} \left| \alpha - \frac{h_n}{k_n} \right| &= \left| \sum_{m=0}^{\infty} \frac{(-1)^m}{k_{m+n} k_{m+n+1}} \right| \\ &= \left| \sum_{m=0}^{\infty} \left(\frac{1}{k_{2m+n} k_{2m+n+1}} - \frac{1}{k_{2m+n+1} k_{2m+n+2}} \right) \right| \\ &= \left| \sum_{m=0}^{\infty} \frac{1}{k_{2m+n+1}} \cdot \frac{k_{2m+n+2} - k_{2m+n}}{k_{2m+n} k_{2m+n+2}} \right|. \end{aligned}$$

Because $\{k_n\}_{n=0}^\infty$ is a strictly increasing sequence, all the terms of the sum are positive, so we may remove the absolute signs to see

$$\left| \alpha - \frac{h_n}{k_n} \right| > \frac{1}{k_{n+1}} \cdot \frac{k_{n+2} - k_n}{k_n k_{n+2}}.$$

Thus, to prove the desired lower bound, we must show $k_{n+1} k_{n+2} < (k_{n+1} + k_n)(k_{n+2} - k_n)$. This rearranges to $k_n^2 < k_n(k_{n+1} + k_{n+2})$, which is true. ■

Remark 1.41. Proposition 1.40 tells us that h_n/k_n will be a “better” rational approximation for α when k_{n+1} is particularly large. For example, $\pi = [3; 7, 15, 1, 292, 1, 1, 1]$, so we would guess that

$$[3; 7, 15, 1] = \frac{355}{113} = 3.14159292035\dots$$

is a particularly good rational approximation of π , and indeed it is. Notably, $[3; 7] = 22/7$ is also a remarkable rational approximation.

As such, we may make the following definition.

Definition 1.42 (infinite continued fraction). Let a_0, a_1, a_2, \dots be integers, where a_1, a_2, \dots are positive. Then we define the *infinite continued fraction*

$$[a_0; a_1, a_2, \dots] := \lim_{n \rightarrow \infty} [a_0; a_1, a_2, \dots, a_n].$$

Example 1.43. We have

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

Solution. By Proposition 1.40, we know that $[1; 1, 1, \dots]$ converges to some real number α . Further,

$$\alpha = 1 + \frac{1}{1 + \frac{1}{1 + \ddots}} = 1 + \frac{1}{\alpha},$$

which rearranges to $\alpha^2 - \alpha - 1 = 0$, so

$$\alpha \in \left\{ \frac{1 \pm \sqrt{5}}{2} \right\}.$$

However, we claim that $\alpha > 0$. With the tools we have, this is somewhat annoying to show, but we remark that Lemma 1.61 makes this relatively easy. Anyway, let $\{h_n/k_n\}_{n=0}^\infty$ denote the continued fraction convergents. Proposition 1.32 implies that $h_0/k_0 = 1/1$ and $h_1/k_1 = 2/1$, so

$$|\alpha - 1| = \left| \alpha - \frac{h_0}{k_0} \right| < \frac{1}{k_0 k_1} = 1,$$

so $\alpha > 0$. Thus, $\alpha = \varphi$. ■

Exercise 1.44. Compute $[2; 2, 2, \dots]$.

The above examples have the amusing feature that $[a_0; a_1, a_2, \dots]$ is irrational. This is not a coincidence. The following result is perhaps our first “Diophantine approximation” result.

Proposition 1.45. Let α be a real number, and let $C > 0$ and $\varepsilon > 0$. Then α is irrational if there is a sequence of rational numbers $\{h_n/k_n\}_{n=0}^\infty$ such that

$$\left| \alpha - \frac{h_n}{k_n} \right| < \frac{C}{k_n^{1+\varepsilon}}$$

for each $n \geq 0$.

Proof. We show the contrapositive. Suppose that $\alpha = p/q$ is rational with $q \geq 1$ and $\gcd(p, q) = 1$, and we show that there are only finitely many rational numbers h/k such that $|\alpha - h/k| < C/k^{1+\varepsilon}$; we may assume that $k \geq 1$ and that $\gcd(h, k) = 1$ in our fractions h/k . Now, for any given k , we note that our inequality rearranges to

$$|h - k\alpha| < \frac{C}{k^\varepsilon},$$

so there are only finitely many integers h in our interval. Thus, it suffices to upper-bound k . Well, plugging in $\alpha = p/q$ and clearing fractions reveals that we want

$$|qh - pk| < \frac{Cq}{k^\varepsilon}.$$

Now, we claim that $k \leq \max \{(Cq)^{1/\varepsilon}, q\}$, which completes the proof. Well, suppose that $k^\varepsilon > Cq$, and we will show $k = q$. Indeed, $qh - pk$ is an integer with magnitude less than 1, so it follows that $qh - pk = 0$, so in fact

$$qh = pk.$$

By the uniqueness of our representation of rational numbers, it follows that $k = q$. Explicitly, $q \mid pk$, but $\gcd(q, p) = 1$, so $q \mid k$. A symmetric argument shows $k \mid q$, so $k, q \geq 1$ establishes $k = q$. ■

Remark 1.46. Proposition 1.45 is fairly surprising result! Approximately speaking, it says that having “too many” good rational approximations of a given real number actually forces the real number to be irrational! We will prove a converse shortly in Corollary 1.53.

Remark 1.47. Here is a way to intuit Proposition 1.45: there is a sense in which rational numbers cannot be “too close to each other” simply because

$$\left| \frac{a}{b} - \frac{c}{d} \right| \geq \frac{1}{|bd|}.$$

Thus, we should not be able to use rational numbers to provide good rational approximations of rational numbers.

Corollary 1.48. Let a_0, a_1, a_2, \dots be integers, where a_1, a_2, \dots are positive. Then $[a_0; a_1, a_2, \dots]$ is irrational.

Proof. Let $\{h_n/k_n\}_{n=0}^\infty$ denote the continued fraction convergents $h_n/k_n := [a_0; a_1, \dots, a_n]$ where $k_n \geq 1$ and $\gcd(h_n, k_n) = 1$. Then Proposition 1.40 establishes that

$$\left| [a_0; a_1, a_2, \dots] - \frac{h_n}{k_n} \right| < \frac{1}{k_n k_{n+1}} < \frac{1}{k_n^2}$$

for each $n \geq 0$, where the last inequality follows because $\{k_n\}_{n=0}^\infty$ is strictly increasing. Proposition 1.45 completes the proof. ■

1.3.2 Building Infinite Continued Fractions

Given an irrational real number $\alpha \in \mathbb{R} \setminus \mathbb{Q}$, we would like to construct a sequence of integers a_0, a_1, a_2, \dots with a_1, a_2, \dots positive and $\alpha = [a_0; a_1, a_2, \dots]$. We did this by hand for φ in Example 1.43, but this is not a general algorithm.

Let’s describe what the algorithm should be. Suppose we could write $\alpha = [a_0; a_1, a_2, \dots]$. Then

$$\alpha = a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \ddots}}$$

forces $a_0 = \lfloor \alpha \rfloor$. From here, define $\alpha_1 := (\alpha - a_0)^{-1}$, and we see

$$\alpha_1 = a_1 + \frac{1}{a_2 + \ddots}.$$

Then we can see that we must have $a_1 = \lfloor \alpha_1 \rfloor$, and we go on to define $\alpha_2 = (\alpha_1 - a_1)^{-1}$ and continue the process. This suggests the following result.

Proposition 1.49. Let $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ be an irrational number. Define the sequence of real numbers $\{\alpha_n\}_{n=0}^\infty$ and integers $\{a_n\}_{n=0}^\infty$ by $\alpha_0 := \alpha$ and

$$a_n := \lfloor \alpha_n \rfloor \quad \text{and} \quad \alpha_{n+1} := \frac{1}{\alpha_n - a_n}$$

Then a_0, a_1, a_2, \dots are integers, and a_1, a_2, \dots are positive, and $\alpha = [a_0; a_1, a_2, \dots]$.

Proof. Quickly, we note that there are no division by zero problems: by construction, the a_n are all integers, and the recursion implies that α_{n+1} is irrational if and only if α_n is irrational, so induction implies that all the α_n are irrational. Next up, we note that $a_n < \alpha_n < a_n + 1$ for each $n \geq 0$ (recall α_n is irrational for each n), so $0 < \alpha_n - a_n < 1$ for each $n \geq 0$, so $a_{n+1} \geq 1$ for each $n \geq 0$, so a_1, a_2, \dots are in fact positive integers.

It remains to show $\alpha = [a_0; a_1, a_2, \dots]$. This is somewhat technical. The main claim is that

$$\alpha \stackrel{?}{=} [a_0; a_1, \dots, a_n, \alpha_{n+1}]$$

for each $n \geq 0$. We show this by induction. For $n = -1$, there is nothing to say because $\alpha = \alpha_0$. For the induction, we write

$$\begin{aligned} \alpha &= [a_0; a_1, \dots, a_n, \alpha_{n+1}] \\ &= [a_0; a_1, \dots, a_n, \lfloor \alpha_{n+1} \rfloor + \{\alpha_{n+1}\}] \\ &= \left[a_0; a_1, \dots, a_n, a_{n+1} + \frac{1}{\alpha_{n+2}} \right] \\ &= [a_0; a_1, \dots, a_n, a_{n+1}, a_{n+2}], \end{aligned}$$

which completes the induction.

We now finish the proof that $\alpha = [a_0; a_1, a_2, \dots]$. For each $n \geq 0$, set $h_n/k_n := [a_0; a_1, \dots, a_n]$ and $h'_{n+1}/k'_{n+1} := [a_0; a_1, a_2, \dots, a_n, \alpha_{n+1}]$ as constructed in Proposition 1.32. Then applying Lemma 1.39 implies

$$\begin{aligned} \alpha - [a_0; a_1, a_2, \dots, a_n] &= [a_0; a_1, \dots, a_n, \alpha_{n+1}] - [a_0; a_1, a_2, \dots, a_n] \\ &= \frac{h_0}{k_0} + \sum_{m=0}^{n-1} \frac{(-1)^m}{k_m k_{m+1}} - \frac{h_0}{k_0} - \sum_{m=0}^{n-1} \frac{(-1)^m}{k_m k_{m+1}} - \frac{(-1)^n}{k_n k'_{n+1}} \\ &= \frac{(-1)^n}{k_n k'_{n+1}}. \end{aligned}$$

Thus,

$$|\alpha - [a_0; a_1, a_2, \dots, a_n]| \leq \frac{1}{k_n^2},$$

where we have used the fact that $k'_{n+1} = \alpha_{n+1} k_n + k_{n-1} \geq k_n$. Sending $n \rightarrow \infty$ makes $k_n \rightarrow \infty$, so we conclude $[a_0; a_1, \dots, a_n] \rightarrow \alpha$ as $n \rightarrow \infty$. ■

Exercise 1.50. Use Proposition 1.49 (and Sage) to compute the first 10 continued fraction coefficients of π .

Remark 1.51. In contrast to Remark 1.26, the continued fraction attached to irrational $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ is unique. The proof is approximately along the lines as the argument at the start of the subsection. Namely, suppose we have integers a_0, a_1, a_2, \dots and b_0, b_1, b_2, \dots with a_1, a_2, \dots and b_1, b_2, \dots positive, and suppose

$$[a_0; a_1, a_2, \dots] = [b_0; b_1, b_2, \dots].$$

We want to show $a_n = b_n$ for all n . Because $[a_0; a_1, a_2, \dots] = a_0 + [a_1; a_2, \dots]^{-1}$, it suffices by induction to show that $a_0 = b_0$. Well, $a_1, b_1 \geq 1$ implies $[a_1; a_2, \dots], [b_1; b_2, \dots] > 1$, so

$$a_0 = \left\lfloor a_0 + \frac{1}{[a_1; a_2, \dots]} \right\rfloor = \lfloor [a_0; a_1, a_2, \dots] \rfloor = \lfloor [b_0; b_1, b_2, \dots] \rfloor = b_0.$$

Proposition 1.49 allows us to make the following terminology.

Definition 1.52 (convergent). Let $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ be an irrational number. By Proposition 1.49, we may find integers a_0, a_1, a_2, \dots where a_1, a_2, \dots are positive and $\alpha = [a_0; a_1, a_2, \dots]$. Then the n th continued fraction convergent of α is $[a_0; a_1, a_2, \dots, a_n]$.

Corollary 1.53. Let $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ be an irrational number. Then there is a sequence of rational numbers $\{h_n/k_n\}_{n=0}^\infty$ such that

$$\left| \alpha - \frac{h_n}{k_n} \right| < \frac{1}{k_n^2}$$

for each $n \geq 0$.

Proof. We use continued fraction convergents. Let $\{h_n/k_n\}_{n=0}^\infty$ be the sequence of continued fraction convergents of α . Then Proposition 1.40 implies

$$\left| \alpha - \frac{h_n}{k_n} \right| < \frac{1}{k_n k_{n+1}}.$$

Because $k_{n+1} > k_n$ by the recursion, the conclusion follows. ■

1.3.3 Quadratic Irrationals

As an intermission, we take a moment to compute the continued fraction of \sqrt{d} where d is a non-square positive integer. Let's start with an example.

Example 1.54. We show that $\sqrt{3} = [1; \overline{1, 2}]$, where the over line indicates periodicity.

Solution. It is possible to solve this by computing $[1; \overline{1, 2}]$ first as in Example 1.43, but we will take a more direct

approach following Proposition 1.49. Indeed, following the algorithm, we compute

$$\begin{aligned}
 \sqrt{3} &= 1 + \left(-1 + \sqrt{3} \right) \\
 &= 1 + \frac{1}{\frac{1+\sqrt{3}}{2}} \\
 &= 1 + \frac{1}{1 + \frac{-1+\sqrt{3}}{2}} \\
 &= 1 + \frac{1}{1 + \frac{1}{1 + \sqrt{3}}} \\
 &= 1 + \frac{1}{1 + \frac{1}{2 + (-1 + \sqrt{3})}}.
 \end{aligned}$$

At this point, we have seen that

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

so $[1, 2] = \frac{1+\sqrt{3}}{2}$ follows, from which the desired result follows by noting $\sqrt{3} = 1 + 1/\left(\frac{1+\sqrt{3}}{2}\right)$ as computed above. ■

As in the above example, it will turn out that the terms α_n from Proposition 1.49 will all take the form

$$\frac{r_n + \sqrt{d}}{s_n}$$

for some positive integers r_n, s_n . We are now ready to state our result.

Proposition 1.55. Fix a non-square positive integer d . Define $a_0 := \lfloor \sqrt{d} \rfloor$ and $r_0 := 0$ and $s_0 := 1$ and $\alpha_0 := \sqrt{d}$ and

$$a_n := \left\lfloor \frac{r_n + a_0}{s_n} \right\rfloor, \quad r_{n+1} := a_n s_n - r_n, \quad \text{and} \quad s_{n+1} := \frac{d - r_{n+1}^2}{s_n}$$

for each $n \geq 0$. Then these are sequences of integers with $0 \leq r_n < \sqrt{d}$ and $1 \leq s_n < 2\sqrt{d}$, and $\sqrt{d} = [a_0; a_1, a_2, \dots]$.

Proof. We forget about the sequences defined in the statement of the proposition, and we will redefine such sequences a different way in the argument which follows. Let $\{a_n\}_{n=0}^\infty$ and $\{\alpha_n\}_{n=0}^\infty$ be as in Proposition 1.49. The main point is to compute the α_n 's. We proceed in steps.

1. We define our sequences. The proof of Proposition 1.49 actually shows that $\sqrt{d} = [a_0; a_1, \dots, a_n; \alpha_{n+1}]$ for any $n \geq 0$, so Proposition 1.32 shows that any $n \geq 0$ has

$$\sqrt{d} = \frac{h_{n-1}\alpha_n + h_{n-2}}{k_{n-1}\alpha_n + k_{n-2}},$$

where $\{h_n/k_n\}_{n=-2}^\infty$ are the continued fraction convergents of \sqrt{d} , and we have taken $(h_{-2}, k_{-2}) = (0, 1)$ and $(h_{-1}, k_{-1}) = (1, 0)$ formally. We can use the above equation to solve α_n as

$$\alpha_n = -\frac{h_{n-2} - k_{n-2}\sqrt{d}}{h_{n-1} - k_{n-1}\sqrt{d}} = -\frac{(h_{n-1}h_{n-2} - dk_{n-1}k_{n-2}) + (-1)^{n-1}\sqrt{d}}{h_{n-1}^2 - dk_{n-1}^2},$$

where we have used Corollary 1.35 in the last equality. As such, we define the integer sequences

$$\begin{aligned} r_n &:= (-1)^n (h_{n-2}h_{n-1} - dk_{n-1}k_n) \\ s_n &:= (-1)^n (h_{n-1}^2 - dk_{n-1}^2) \end{aligned}$$

so that

$$\alpha_n = \frac{r_n + \sqrt{d}}{s_n}.$$

The reason we have chosen to define r_\bullet and s_\bullet this way is that we know that they are integer sequences "for free."

2. We show the recursions for r_\bullet and s_\bullet . To begin, $\alpha_0 = \sqrt{d}$ forces $(r_0, s_0) = (0, 1)$. Further, the remaining recursion comes from the recursion of Proposition 1.49: write

$$\frac{r_{n+1} + \sqrt{d}}{s_{n+1}} = \alpha_{n+1} = \frac{1}{\alpha_n - a_n} = \frac{1}{\frac{r_n + \sqrt{d}}{s_n} - a_n} = \frac{s_n}{-(a_n s_n - r_n) + \sqrt{d}} = \frac{(a_n s_n - r_n) + \sqrt{d}}{\frac{d - (a_n s_n - r_n)^2}{s_n}}.$$

Thus, comparing coefficients, we see $r_{n+1} = a_n s_n - r_n$ and $s_{n+1} = \frac{1}{s_n} (d - r_{n+1}^2)$, as needed.

3. We show the inequalities on r_\bullet and s_\bullet . Quickly, note that $h_{n-1}^2 - dk_{n-1}^2 > 0$ if and only if $h_{n-1}/k_{n-1} > \sqrt{d}$ if and only if $n-1$ is odd by Lemma 1.61, which is equivalent to n being even, meaning that $s_n > 0$ always. The recursion on s_\bullet then forces $r_{n+1} < \sqrt{d}$ always, and combined with $r_0 = 0$, we see that $r_n < \sqrt{d}$ always. For the other inequalities, we show

$$0 < \frac{\sqrt{d} - r_n}{s_n} < 1$$

for $n \geq 1$. For $n = 1$, we see $r_1 = a_0$ and $s_1 = d - a_0^2$, so the given quantity is $1/(\sqrt{d} + a_0) < 1$. As for the inductive step, by replacing \sqrt{d} with $-\sqrt{d}$ in the computation of the previous step, we see

$$\frac{r_{n+1} - \sqrt{d}}{s_{n+1}} = \frac{1}{\frac{r_n - \sqrt{d}}{s_n} - a_n},$$

so

$$\frac{\sqrt{d} - r_{n+1}}{s_{n+1}} = \frac{1}{\frac{\sqrt{d} - r_n}{s_n} + a_n}$$

is positive and less than 1 because $a_n \geq 1$ and the inductive hypothesis.

Now, $\alpha_n = \frac{r_n + \sqrt{d}}{s_n} > 1$ by the proof of Proposition 1.49, so adding and subtracting yields $\frac{2r_n}{s_n} > 0$ and $\frac{2\sqrt{d}}{s_n} > 1$, so $r_n > 0$ and $s_n < 2\sqrt{d}$ for $n > 0$, thus completing the step.

4. We show the recursion on a_\bullet . This merely requires writing

$$a_n = \lfloor \alpha_n \rfloor = \left\lfloor \frac{r_n + \sqrt{d}}{s_n} \right\rfloor = \left\lfloor \frac{r_n + a_0}{s_n} \right\rfloor,$$

so we are done. ■

Corollary 1.56. Fix a non-square positive integer d , and write $\sqrt{d} = [a_0; a_1, a_2, \dots]$. Then there exists an integer $N \leq 2d$ and positive integer $p \leq 2d$ such that $a_{n+p} = a_n$ for each $n \geq N$.

Proof. We use Proposition 1.55. Among the d pairs $(r_0, s_0), (r_1, s_1), \dots, (r_{2d}, s_{2d})$, we must repeat some ordered pair of integers; say $(r_N, s_N) = (r_{N+p}, s_{N+p})$ for some integer $N \leq 2d$ and positive integer $p < 2d$. Then the recursion

$$(r_{n+1}, s_{n+1}) = \left(\left\lfloor \frac{r_n + a_0}{s_n} \right\rfloor s_n - r_n, \frac{d - r_n^2}{s_n} \right)$$

forces $(r_n, s_n) = (r_{n+p}, s_{n+p})$ for each $n \geq N$, so $a_{n+p} = \left\lfloor \frac{r_{n+p} + a_0}{s_{n+p}} \right\rfloor = \left\lfloor \frac{r_n + a_0}{s_n} \right\rfloor = a_n$ follows. ■

Remark 1.57. Examining the above proofs, we actually see that

$$\begin{aligned} r_n &:= (-1)^n (h_{n-2}h_{n-1} - dk_{n-1}k_n) \\ s_n &:= (-1)^n (h_{n-1}^2 - dk_{n-1}^2) \end{aligned}$$

enjoys the same periodicity as Corollary 1.56.

1.3.4 Convergents Are Good Rational Approximations

As before, let a_0, a_1, a_2, \dots be integers, where a_1, a_2, \dots are positive, and let $\{h_n/k_n\}_{n=0}^\infty$ denote the continued fraction convergents $h_n/k_n := [a_0; a_1, \dots, a_n]$ where $k_n \geq 1$ and $\gcd(h_n, k_n) = 1$. Proposition 1.40 immediately implies that

$$\left| \alpha - \frac{h_n}{k_n} \right| \leq \frac{1}{k_n^2},$$

but we can improve this result somewhat. The goal of the present section is to show that there are infinitely many n for which

$$\left| \alpha - \frac{h_n}{k_n} \right| \leq \frac{1}{\sqrt{5}k_n^2},$$

and the following example explains that the constant $\sqrt{5}$ is the best possible.

Example 1.58. Let $\varphi = \frac{1+\sqrt{5}}{2} = [1; 1, 1, \dots]$ as in Example 1.43. By Example 1.34, the n th continued fraction convergent is F_{n+2}/F_{n+1} . For any $c > \sqrt{5}$, we have

$$\left| \varphi - \frac{F_{n+2}}{F_{n+1}} \right| < \frac{1}{cF_{n+1}^2}$$

for only finitely many n .

Solution. Set $\bar{\varphi} := \frac{1-\sqrt{5}}{2}$, which is the negative solution of $x^2 = x + 1$; note $\varphi + \bar{\varphi} = 1$ and $\varphi\bar{\varphi} = -1$. An induction n proves Binet's formula

$$F_n = \frac{\varphi^n - \bar{\varphi}^n}{\sqrt{5}}.$$

Indeed, the above equality holds at $n = 0$ and $n = 1$ by a direct computation, and taking a linear combination of the relations $\varphi^{n+2} = \varphi^{n+1} + \varphi^n$ and $\bar{\varphi}^{n+2} = \bar{\varphi}^{n+1} + \bar{\varphi}^n$ proves the inductive step.

We now carefully study the error. For any $n \geq 0$, we see

$$\begin{aligned} 5(\varphi F_{n+1}^2 - F_{n+2}F_{n+1}) &= \varphi(\varphi^{n+1} - \bar{\varphi}^{n+1})^2 - (\varphi^{n+2} - \bar{\varphi}^{n+2})(\varphi^{n+1} - \bar{\varphi}^{n+1}) \\ &= \varphi(\varphi^{2n+2} + \bar{\varphi}^{2n+2} - 2(\varphi\bar{\varphi})^{n+1}) - (\varphi^{2n+3} + \bar{\varphi}^{2n+3} - (\varphi\bar{\varphi})^{n+1}(\varphi + \bar{\varphi})) \\ &= (-1)^n(2\varphi - 1) + \bar{\varphi}^{2n+2}(\varphi - \bar{\varphi}) \\ &= (-1)^n\sqrt{5} + \bar{\varphi}^{2n+2}\sqrt{5}. \end{aligned}$$

Thus,

$$cF_{n+1}^2 \left| \varphi - \frac{F_{n+2}}{F_{n+1}} \right| = \frac{c}{\sqrt{5}} |(-1)^n + \bar{\varphi}^{2n+2}|. \quad (1.1)$$

As $n \rightarrow \infty$, we see $\bar{\varphi}^{2n+2} \rightarrow 0$, so the error above approaches $c/\sqrt{5} > 1$. Thus, only finitely many n have the above quantity less than 1, which is what we wanted. ■

Remark 1.59. Carefully tracking through Example 1.58 tells us that

$$\left| \varphi - \frac{F_{n+2}}{F_{n+1}} \right| < \frac{1}{\sqrt{5}F_{n+1}^2}$$

exactly for the even n . Indeed, this follows from (1.1) upon noting $-\bar{\varphi}^{2n+2} < 0$. Compare this result with the statement and proof of Theorem 1.63.

Exercise 1.60. Set $\alpha := \sqrt{2}$, and let $\{h_n/k_n\}_{n=0}^\infty$ be the continued fraction convergents of α . Find the largest real number $c > 0$ for which there exist infinitely many integers $n \geq 0$ such that

$$\left| \alpha - \frac{h_n}{k_n} \right| < \frac{1}{ck_n^2}.$$

As should be somewhat evident by the $\sqrt{5}$ in our bounds and in the above proof, the arguments here are going to be somewhat ad-hoc. The following result starts us off.

Lemma 1.61. Let $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ be irrational, and let $\{h_n/k_n\}_{n=0}^\infty$ be the sequence of continued fraction convergents of α . For any $n \geq 0$, we have

$$\frac{h_{2n}}{k_{2n}} < \frac{h_{2n+2}}{k_{2n+2}} < \frac{h_{2n+3}}{k_{2n+3}} < \frac{h_{2n+1}}{k_{2n+1}}.$$

Proof. Applying Lemma 1.39, we are trying to show

$$\frac{h_{2n}}{k_{2n}} \stackrel{?}{<} \frac{h_{2n}}{k_{2n}} + \frac{1}{k_{2n}k_{2n+1}} - \frac{1}{k_{2n+1}k_{2n+2}} \stackrel{?}{<} \frac{h_{2n}}{k_{2n}} + \frac{1}{k_{2n}k_{2n+1}} - \frac{1}{k_{2n+1}k_{2n+2}} + \frac{1}{k_{2n+2}k_{2n+3}} \stackrel{?}{<} \frac{h_{2n}}{k_{2n}} + \frac{1}{k_{2n}k_{2n+1}}.$$

Simplifying, we want to show

$$0 \stackrel{?}{<} \frac{1}{k_{2n}k_{2n+1}} - \frac{1}{k_{2n+1}k_{2n+2}} \stackrel{?}{<} \frac{1}{k_{2n}k_{2n+1}} - \frac{1}{k_{2n+1}k_{2n+2}} + \frac{1}{k_{2n+2}k_{2n+3}} \stackrel{?}{<} \frac{1}{k_{2n}k_{2n+1}}.$$

The leftmost inequality is equivalent to $k_{2n} < k_{2n+2}$, which is true. The middle inequality is equivalent to $0 < 1/(k_{2n+2}k_{2n+3})$, which is true. Lastly, the rightmost inequality is equivalent to $k_{2n+1} < k_{2n+3}$, which is true. ■

Proposition 1.62. Let $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ be irrational, and let $\{h_n/k_n\}_{n=0}^\infty$ be the sequence of continued fraction convergents of α . For any $m \geq 0$, there exists $n \in \{2m, 2m+1\}$ such that

$$\left| \alpha - \frac{h_n}{k_n} \right| < \frac{1}{2k_n^2}.$$

Proof. The point is that one of h_{2m}/k_{2m} or h_{2m+1}/k_{2m+1} is going to be “closer” to α . By Lemma 1.61, we see that $\{h_{2m}/k_{2m}\}_{m=0}^{\infty}$ is a strictly ascending sequence of rational numbers, which converges to α by definition of α . Analogously, $\{h_{2m+1}/k_{2m+1}\}_{m=0}^{\infty}$ is a strictly descending sequence of rational numbers which also converges to α . Thus,

$$\frac{h_{2m}}{k_{2m}} < \alpha < \frac{h_{2m+1}}{k_{2m+1}}.$$

By Lemma 1.39, the length of this interval is $1/(k_{2m}k_{2m+1})$.

Now, suppose for contradiction that

$$\left| \alpha - \frac{h_n}{k_n} \right| \geq \frac{1}{2k_n^2}$$

for $n \in \{2m, 2m+1\}$. Then we must have

$$\frac{h_{2m}}{k_{2m}} + \frac{1}{2k_{2m}^2} \leq \alpha \leq \frac{h_{2m+1}}{k_{2m+1}} - \frac{1}{2k_{2m+1}^2}.$$

This rearranges to

$$\frac{1}{2k_{2m}^2} + \frac{1}{2k_{2m+1}^2} \leq \frac{1}{k_{2m}k_{2m+1}}$$

by Lemma 1.39, but this is equivalent to $(k_{2m} - k_{2m+1})^2 \leq 0$, or $k_{2m} = k_{2m+1}$. This is a contradiction because the sequence $\{k_n\}_{n=0}^{\infty}$ is strictly increasing. ■

With a little more care in the last half of the argument, we can achieve the desired result.

Theorem 1.63 (Hurwitz). Let $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ be irrational, and let $\{h_n/k_n\}_{n=0}^{\infty}$ be the sequence of continued fraction convergents of α . For any $m \geq 0$, there exists $n \in \{3m, 3m+1, 3m+2\}$ such that

$$\left| \alpha - \frac{h_n}{k_n} \right| < \frac{1}{\sqrt{5}k_n^2}.$$

Proof. The proof is along the same lines as Proposition 1.62. Without loss of generality, we work with even m in order to make our inequalities better-behaved; the argument for odd m is analogous but requires reversing a few inequalities. Anyway, if m is even, Lemma 1.61 implies

$$\frac{h_{3m}}{k_{3m}} < \frac{h_{3m+2}}{k_{3m+2}} < \alpha < \frac{h_{3m+1}}{k_{3m+1}}.$$

(The location of α adjusts in the case where m is odd.) Now, suppose for the sake of contradiction that

$$\left| \alpha - \frac{h_n}{k_n} \right| \geq \frac{1}{\sqrt{5}k_n^2}$$

for each $n \in \{3m, 3m+1, 3m+2\}$. Removing the absolute values, we receive the inequalities

$$\frac{h_{3m}}{k_{3m}} + \frac{1}{\sqrt{5}k_{3m}^2} \leq \alpha, \quad \alpha \leq \frac{h_{3m+1}}{k_{3m+1}} - \frac{1}{\sqrt{5}k_{3m+1}^2}, \quad \text{and} \quad \frac{h_{3m+2}}{k_{3m+2}} + \frac{1}{\sqrt{5}k_{3m+2}^2} \leq \alpha,$$

which imply

$$\frac{h_{3m}}{k_{3m}} + \frac{1}{\sqrt{5}k_{3m}^2} \leq \frac{h_{3m+1}}{k_{3m+1}} - \frac{1}{\sqrt{5}k_{3m+1}^2}, \quad \text{and} \quad \frac{h_{3m+2}}{k_{3m+2}} + \frac{1}{\sqrt{5}k_{3m+2}^2} \leq \frac{h_{3m+1}}{k_{3m+1}} - \frac{1}{\sqrt{5}k_{3m+1}^2}.$$

By Lemma 1.39, these rearrange into

$$\frac{1}{k_{3m}^2} + \frac{1}{k_{3m+1}^2} \leq \frac{\sqrt{5}}{k_{3m}k_{3m+1}}, \quad \text{and} \quad \frac{1}{k_{3m+1}^2} + \frac{1}{k_{3m+2}^2} \leq \frac{\sqrt{5}}{k_{3m+1}k_{3m+2}}.$$

By Proposition 1.32, we see that $k_{3m} + k_{3m+1} \leq k_{3m+2}$, so our inequalities read

$$\frac{1}{k_{3m}^2} + \frac{1}{k_{3m+1}^2} \leq \frac{\sqrt{5}}{k_{3m}k_{3m+1}}, \quad \text{and} \quad \frac{1}{k_{3m+1}^2} + \frac{1}{(k_{3m} + k_{3m+1})^2} \leq \frac{\sqrt{5}}{k_{3m+1}(k_{3m} + k_{3m+1})}.$$

Now, we set $q := k_{3m+1}/k_{3m}$ to homogenize the inequalities. This gives

$$q^2 + 1 \leq \sqrt{5}q, \quad \text{and} \quad (q + 1)^2 + 1 \leq \sqrt{5}(q + 1).$$

In other words, we are asking for $\{q, q+1\} \subseteq \{x \in \mathbb{R} : x^2 + 1 \leq \sqrt{5}x\}$. To solve for q , we note $x^2 - \sqrt{5}x + 1 = 0$ exactly when $x = \frac{\sqrt{5} \pm 1}{2}$, so $\{x \in \mathbb{R} : x^2 + 1 \leq \sqrt{5}x\}$ is the closed interval from $\frac{\sqrt{5}-1}{2}$ up to $\frac{\sqrt{5}+1}{2}$. Thus, we must have $q = \frac{\sqrt{5}-1}{2}$, which is a contradiction because q is rational while $\frac{\sqrt{5}-1}{2}$ is irrational! ■

1.3.5 Convergents Are Best Rational Approximations

Now that we are somewhat acquainted with what it means to be a “good” rational approximation, we are ready to state and prove our main result on continued fractions. It is a converse to Proposition 1.62. Our exposition in this subsection roughly follows [HW75, Theorem 184].

Theorem 1.64. Let $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ be irrational, and let $\{h_n/k_n\}_{n=0}^\infty$ be the sequence of continued fraction convergents of α . Given a rational number h/k with $\gcd(h, k) = 1$ and $k \geq 1$, if

$$\left| \alpha - \frac{h}{k} \right| < \frac{1}{2k^2},$$

then $(h, k) = (h_n, k_n)$ for some n .

Approximately speaking, Theorem 1.64 tells us that the best rational approximations of a real number are all continued fraction convergents.

Proof of Theorem 1.64. We use Remark 1.26 to write

$$\frac{h}{k} = [a_0; a_1, a_2, \dots, a_n]$$

with a_n with parity chosen so that n is even if and only if $\alpha > h/k$. (This is what we expect from Lemma 1.61.) Then let $\{p_m/q_m\}_{m=0}^n$ be the continued fraction convergents; for example, $(h, k) = (p_n, q_n)$.

The main idea is to show that the continued fraction expansion of α begins $[a_0; a_1, a_2, \dots, a_n, \dots]$. To realize this, we must continue the continued fraction. Well, we know that we can certainly find some $\beta \in \mathbb{R}$ such that

$$\alpha = \frac{p_n\beta + p_{n-1}}{q_n\beta + q_{n-1}}$$

by rearranging. (Explicitly, we need to know that $\alpha k - h \neq 0$ to set $\beta := (h' - \alpha k')/(\alpha k - h)$, which is true because α is irrational.) The main claim is that $\beta > 1$. Well, comparing with our error, we see

$$\alpha - \frac{h}{k} = \frac{p_n\beta + p_{n-1}}{q_n\beta + q_{n-1}} - \frac{p_n}{q_n} = \frac{p_{n-1}q_n - p_nq_{n-1}}{(q_n\beta + q_{n-1})q_n} = \frac{(-1)^n}{(q_n\beta + q_{n-1})q_n},$$

where we applied Corollary 1.35 in the last equality. We arranged the parity n so that the left-hand side is positive if and only if $(-1)^n = 1$, so we may now write

$$1 > 2p_n^2 \left| \alpha - \frac{p_n}{q_n} \right| = \frac{2p_n}{p_n\beta + p_{n-1}},$$

so $\beta > 2 - p_{n-1}/p_n$, which is bigger than 1 because $p_{n-1} < p_n$.

We now convert $\beta > 1$ into the result. Well, Proposition 1.49 allows us to write

$$\beta = [a_{n+1}; a_{n+2}, a_{n+3}, \dots]$$

for integers $a_{n+1}, a_{n+2}, a_{n+3}, \dots$ with a_{n+2}, a_{n+3}, \dots positive. In fact, $a_{n+1} = \lfloor \beta \rfloor \geq 1$ is positive by construction; here is where we used $\beta > 1$. We conclude that

$$\alpha = \frac{p_n \beta + p_{n-1}}{q_n \beta + q_{n-1}} = [a_0; a_1, \dots, a_n, \beta] = [a_0; a_1, \dots, a_n, a_{n+1}, a_{n+2}, \dots].$$

By the uniqueness of the continued fraction (see Remark 1.51), we conclude that $(p_m, q_m) = (h_m, k_m)$ for $0 \leq m \leq n$, which completes the proof upon setting $m = n$. ■

1.3.6 Problems

Do at least ten points worth of the following exercises.

Problem 1.3.1 (2 points). Work Exercise 1.44.

Problem 1.3.2 (3 points). Work Exercise 1.60.

Problem 1.3.3 (3 points). Let a_0, a_1, a_2, \dots be integers with a_1, a_2, \dots positive. Suppose that there exists an integer m such that $a_n = a_{n+m}$ for all n . Show that $[a_0; a_1, a_2, \dots]$ is the root of a polynomial with integer coefficients and of degree two.

Problem 1.3.4 (4 points). Find an irrational number $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ and integers h and k such that

$$\left| \alpha - \frac{h}{k} \right| < \frac{1}{k^2},$$

but h/k is not a continued fraction convergent of α .

Problem 1.3.5 (5 points). Write (and submit) a Python program which takes as input an irrational number $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ and an index n and then outputs the n th coefficient a_n of the corresponding continued fraction $[a_0; a_1, a_2, \dots]$ equal to α .

Problem 1.3.6 (8 points). Let $\alpha \in \mathbb{R}$ be irrational and $[a_0; a_1, a_2, \dots]$ its continued fraction expansion. Fix N sufficiently large. Suppose that among the first $1000N$ digits of the decimal expansion of α , the last $999N$ of them are all zeroes or all nines. Then there exists some $n \leq 5N$ so that $a_n > 10^{100N}$.

Problem 1.3.7 (2 points). Use Problem 1.3.6 to conclude that for any sufficiently large N , the last $999N$ digits of the first $1000N$ decimal digits in the decimal expansion of $\sqrt{5}$ cannot be all zeroes or all nines.

1.4 Diophantine Approximation

Now that we have some experience with finding good rational approximations to real numbers, we are able to more firmly step foot into the field of Diophantine approximation. The content in this section is more intensive than in previous sections because it is essentially topics in Diophantine approximation.

1.4.1 Irrationality Measure

Fix an irrational number $\alpha \in \mathbb{R} \setminus \mathbb{Q}$. From one perspective, the arc of the previous section was to go from knowing that there are infinitely rational numbers h/k such that

$$\left| \alpha - \frac{h}{k} \right| < \frac{1}{k}$$

to knowing that there are infinitely many rational numbers h/k such that

$$\left| \alpha - \frac{h}{k} \right| < \frac{1}{k^2}.$$

This is an amazing improvement: going from k to k^2 is a full exponent! But then we spent a lot of time improving the above result into

$$\left| \alpha - \frac{h}{k} \right| < \frac{1}{\sqrt{5}k^2},$$

which feels less significant because we are only improving by a constant. Of course, Example 1.58 established that we cannot do better than this in general, but for some real numbers, it will be possible. With this in mind, we take the following definition.

Definition 1.65 (irrationality measure). Fix a real number $\alpha \in \mathbb{R}$. Then the *irrationality measure* $\mu(\alpha)$ of α is the least upper bound on the set of real numbers r such that

$$\left| \alpha - \frac{h}{k} \right| < \frac{1}{k^r}$$

for infinitely many rational numbers h/k with $k > 0$. Note that we allow $\mu(\alpha) = \infty$.

Remark 1.66. Note that there always is some real number r such that $\left| \alpha - \frac{h}{k} \right| < \frac{1}{k^r}$ for infinitely many rational numbers h/k , which makes the above definition make sense. Indeed, we may take $r = 1$. To see this, for any positive integer k , set $h := \lfloor k\alpha \rfloor$ as in the previous section, so we find

$$\left| \alpha - \frac{h}{k} \right| = \frac{|k\alpha - \lfloor k\alpha \rfloor|}{k} < \frac{1}{k}.$$

So there are indeed infinitely many rational numbers h/k such that $\left| \alpha - \frac{h}{k} \right| < \frac{1}{k}$ as we let k vary.

Here are some early examples.

Example 1.67. Let α be a rational number. Then $\mu(\alpha) = 1$.

Solution. Remark 1.66 establishes $\mu(\alpha) \geq 1$. Further, for any $r > 1$, there are only finitely many rational numbers h/k such that

$$\left| \alpha - \frac{h}{k} \right| < \frac{1}{k^r}$$

by Proposition 1.45. Thus, $\mu(\alpha) \leq 1$, so the result follows. ■

Lemma 1.68. Let $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ be an irrational number. Then $\mu(\alpha) \geq 2$.

Proof. This follows directly from Corollary 1.53 upon unwinding. ■

Example 1.69. We have $\mu(\varphi) = 2$, where $\varphi = \frac{1+\sqrt{5}}{2}$.

Solution. Lemma 1.68 tells us that $\mu(\varphi) \geq 2$, so we just need to show that $\mu(\varphi) \leq 2$. It suffices to show that, for any $\varepsilon > 0$, there are only finitely many rational numbers h/k such that

$$\left| \varphi - \frac{h}{k} \right| < \frac{1}{k^{2+2\varepsilon}}.$$

Well, for sufficiently large k , we have $2k^{2+\varepsilon} < k^{2+2\varepsilon}$, so it is enough to show that there are finitely many rational numbers h/k such that

$$\left| \varphi - \frac{h}{k} \right| < \frac{1}{2k^{2+\varepsilon}}.$$

By Theorem 1.64, all such rational numbers h/k are continued fraction convergents. Thus, we take a moment to recall that $\{F_{n+2}/F_{n+1}\}_{n=0}^{\infty}$ are the continued fraction convergents of φ by Example 1.34, so it is enough to show that there are finitely many nonnegative integers n such that

$$\left| \varphi - \frac{F_{n+2}}{F_{n+1}} \right| < \frac{1}{2F_{n+1}^{2+\varepsilon}}.$$

However, Proposition 1.40 tells us that any nonnegative integer n has

$$\frac{1}{F_{n+1}(F_{n+1} + F_{n+2})} < \left| \varphi - \frac{F_{n+2}}{F_{n+1}} \right|,$$

so rearranging implies that it is enough to show there are only finitely many n with

$$2F_{n+1}^{\varepsilon} < \frac{F_{n+1} + F_{n+2}}{F_{n+1}} = 1 + \frac{F_{n+2}}{F_{n+1}}.$$

However, $F_{n+2}/F_{n+1} \rightarrow \varphi$ as $n \rightarrow \infty$, so the right-hand side is bounded while the left-hand side is not, so indeed there can be only finitely many n satisfying the above inequality. ■

Exercise 1.70. Show that $\mu(\sqrt{2}) = 2$.

Remark 1.71. It is not too hard to show that the continued fraction expansion for any quadratic irrational number α is eventually periodic, so the arguments of the previous two examples show that $\mu(\alpha) = 2$.

Example 1.72 (Liouville). The real number

$$L := \sum_{k=0}^{\infty} \frac{1}{2^{k!}}$$

has $\mu(L) = +\infty$.

Proof. Quickly, note that the series converges because it is bounded above by $\sum_{k=0}^{\infty} 1/2^k = 2$. Now, for each natural n , define

$$L_n := \sum_{k=0}^n \frac{1}{2^{k!}}$$

to be the n th partial sum of L . Then L_n is a rational number with denominator $2^{n!}$, but

$$|L - L_n| = \sum_{k=n+1}^{\infty} \frac{1}{2^{k!}} < \sum_{k=(n+1)!}^{\infty} \frac{1}{2^k} = \frac{1}{2^{(n+1)!-1}}. \quad (1.2)$$

We are now ready to claim that $\mu(L) > r$ for any real number r . Indeed, for any real number r , we claim that there are infinitely many rational numbers h/k such that $|\alpha - h/k| < 1/k^r$. In fact, we claim that there are infinitely many n such that

$$|\alpha - L_n| < \frac{1}{2^{rn!}}.$$

Indeed, (1.2) implies that it is enough to show that

$$\frac{1}{2^{(n+1)!-1}} < \frac{1}{2^{rn!}}$$

for n sufficiently large, which is equivalent to $rn! < (n+1)! - 1$ for n sufficiently large, which is equivalent to $r < n + 1 - 1/n!$ for n sufficiently large, which is true. ■

Before continuing, we should note that essentially all real numbers α have irrationality measure 2. In particular, Example 1.69 is typical, and Example 1.72 is highly remarkable.

Proposition 1.73. Almost all real numbers $\alpha \in \mathbb{R}$ have $\mu(\alpha) = 2$. In other words, for any $\varepsilon > 0$, there is a countable collection of bounded intervals $\{(a_n, b_n)\}_{n=0}^{\infty}$ containing all $\alpha \in \mathbb{R}$ with $\mu(\alpha) = 2$ but

$$\sum_{n=0}^{\infty} (b_n - a_n) < \varepsilon.$$

Proof. Let S be the set of all $\alpha \in \mathbb{R}$ with $\mu(\alpha) \neq 2$. For brevity, let a subset $N \subseteq \mathbb{R}$ be a “null set” if and only if, for all $\varepsilon > 0$, there is a countable collection of bounded intervals $\{(a_n, b_n)\}_{n=0}^{\infty}$ containing N such that

$$\sum_{n=0}^{\infty} (b_n - a_n) < \varepsilon.$$

For example, we see that the union of countably many null sets is a null set by combining the countable collections $\{(a_n, b_n)\}_{n=0}^{\infty}$ together. Additionally, a point $\{x\}$ is a null set because x is covered by $(x - \varepsilon/2, x + \varepsilon/2)$ for any $\varepsilon > 0$. It follows from the previous two sentences that \mathbb{Q} is a null set (it’s a countable union of points), and thus it is enough to show that $S \setminus \mathbb{Q}$ is a null set.

Now, by Lemma 1.68, we see that $\alpha \in S \setminus \mathbb{Q}$ must have $\mu(\alpha) > 2$. For any real number $\varepsilon > 0$, we claim that

$$S_{\varepsilon} := \{\alpha \in \mathbb{R} : \mu(\alpha) > 2 + \varepsilon\}$$

is a null set. By taking the union of $S = S_1 \cup S_{1/2} \cup S_{1/3} \cup \dots$, it will follow that S is a null set. By taking another countable union, it is enough to show that $S_{\varepsilon, M} := S_{\varepsilon} \cap [-M, M]$ is a null set for any $M > 0$. Well, any $\alpha \in S_{\varepsilon}$ has infinitely many rational numbers h/k such that

$$\left| \alpha - \frac{h}{k} \right| < \frac{1}{k^{2+\varepsilon}}.$$

In particular, $\alpha \in S_{\varepsilon, M}$ is contained in the set

$$S_{\varepsilon, M, K} := \left\{ \alpha \in [-M, M] : \left| \alpha - \frac{h}{k} \right| < \frac{1}{k^{2+\varepsilon}} \text{ for some } h \in \mathbb{Z} \text{ and } k \in \mathbb{Z} \cap [K, \infty) \right\}$$

for any $K > 1$. However, the above set is a countable union of small intervals: for each integer $k > K$, the set of relevant α have $2Mk + 1$ options for h , and then each h has an interval of length $2/k^{2+\varepsilon}$ around it. The point is that $S_{\varepsilon, M, K}$ is covered by a countable union of intervals whose lengths total

$$\sum_{k>K} (2Mk + 1) \cdot \frac{2}{k^{2+\varepsilon}} < \sum_{k>K} 3Mk \cdot \frac{2}{k^{2+\varepsilon}} = 6M \sum_{k>K} \frac{1}{k^{1+\varepsilon}}.$$

However, the series $\sum_{k=1}^{\infty} 1/k^{1+\varepsilon}$ converges, so the above length must go to zero as $K \rightarrow \infty$. Thus, for any $\delta > 0$, we can find K large enough so that $S_{\varepsilon, M, K}$ is covered by a countable union of intervals whose lengths sum to less than δ ; this means that $S_{\varepsilon, M}$ is a null set, completing the proof. ■

1.4.2 Irrationality Measure via Continued Fractions

Example 1.69 has reminded us of the important fact that continued fraction convergents provide the best rational approximations. Thus, we might expect the irrationality measure to be controlled by the continued fraction convergents, which is indeed the case.

Lemma 1.74. Let $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ be an irrational number with continued fraction convergents $\{h_n/k_n\}_{n=0}^{\infty}$. Then

$$\mu(\alpha) = \limsup_{n \rightarrow \infty} \frac{-\log |\alpha - h_n/k_n|}{\log k_n}.$$

Proof. Let the lim sup be L . Quickly, recall that Proposition 1.40 implies

$$\frac{-\log |\alpha - p_n/q_n|}{\log q_n} \geq \frac{-\log (1/q_n^2)}{\log q_n} = 2$$

for all n , so $L \geq 2$ has some lower bound.

For a given real number r , we claim that $\mu(\alpha) > r$ if and only if $L > r$, which complete the proof. Note that we may assume $r \geq 2$ because $\mu(\alpha) \geq 2$ by Lemma 1.68 and $L \geq 2$ already. Well, $\mu(\alpha) > r$ is equivalent to having some $\varepsilon > 0$ such that there are infinitely many rational numbers h/k such that

$$\left| \alpha - \frac{h}{k} \right| < \frac{1}{k^{r+\varepsilon}}.$$

Because $r \geq 2$, we see that $r + \varepsilon > 2$. Additionally, for k large enough, we have $2k^2 < 2k^{r+\varepsilon/2} < k^{r+\varepsilon}$, so all sufficiently large k must have h/k a continued fraction convergent. As such, this is equivalent to having infinitely many nonnegative integers n such that

$$\left| \alpha - \frac{h_n}{k_n} \right| < \frac{1}{k_n^{r+\varepsilon}},$$

or

$$\frac{-\log |\alpha - h_n/k_n|}{\log k_n} > r + \varepsilon.$$

This is now equivalent to $L \geq r + \varepsilon$ for our $\varepsilon > 0$, which is equivalent to $L > r$, completing the proof. ■

Proposition 1.75. Let $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ be an irrational number with continued fraction $[a_0; a_1, a_2, \dots]$ and convergents $\{h_n/k_n\}_{n=0}^{\infty}$. Then

$$\mu(\alpha) = 1 + \limsup_{n \rightarrow \infty} \frac{\log k_{n+1}}{\log k_n} = 2 + \limsup_{n \rightarrow \infty} \frac{\log a_{n+1}}{\log k_n}.$$

Proof. We show the equalities separately.

- The left equality follows from Lemma 1.74. To see this, note that Proposition 1.40 implies

$$\frac{1}{2k_n k_{n+1} + 1} < \frac{1}{k_n(k_n + k_{n+1})} < \left| \alpha - \frac{h_n}{k_n} \right| < \frac{1}{k_n k_{n+1}}$$

for any nonnegative integer n . Thus, Lemma 1.74 implies

$$\limsup_{n \rightarrow \infty} \frac{\log 2k_n k_{n+1}}{\log k_n} \geq \mu(\alpha) \geq \limsup_{n \rightarrow \infty} \frac{\log k_n k_{n+1}}{\log k_n},$$

which is equivalent to

$$1 + \limsup_{n \rightarrow \infty} \left(\frac{\log k_{n+1}}{\log k_n} + \frac{\log 2}{\log k_n} \right) \geq \mu(\alpha) \geq 1 + \limsup_{n \rightarrow \infty} \frac{\log k_{n+1}}{\log k_n},$$

For any $\varepsilon > 0$, there is N big enough so that $\log 2 / \log k_n < \varepsilon$ for $n > N$, meaning

$$1 + \varepsilon + \limsup_{n \rightarrow \infty} \frac{\log k_{n+1}}{\log k_n} \geq \mu(\alpha) \geq 1 + \limsup_{n \rightarrow \infty} \frac{\log k_{n+1}}{\log k_n}.$$

Sending $\varepsilon \rightarrow 0^+$ completes the proof.

- The right equality follows from Proposition 1.32. To see this, recall from Proposition 1.32 that

$$k_{n+1} = a_{n+1}k_n + k_{n-1} = a_{n+1}k_n \left(1 + \frac{k_{n-1}}{a_{n+1}k_n} \right)$$

for $n \geq 1$, so

$$\limsup_{n \rightarrow \infty} \frac{\log k_{n+1}}{\log k_n} = 1 + \limsup_{n \rightarrow \infty} \left(\frac{\log a_{n+1}}{\log k_n} + \frac{\log \left(1 + \frac{k_{n-1}}{a_{n+1}k_n} \right)}{\log k_n} \right).$$

Notably, $0 \leq k_{n-1}/(a_{n+1}k_n) \leq 1$ for all $n \geq 1$, so we conclude that

$$\limsup_{n \rightarrow \infty} \frac{\log a_{n+1}}{\log k_n} \leq \mu(\alpha) - 2 \leq \limsup_{n \rightarrow \infty} \left(\frac{\log a_{n+1}}{\log k_n} + \frac{\log 2}{\log k_n} \right).$$

Now, for any $\varepsilon > 0$, we can find N so that $\log 2 / \log k_n < \varepsilon$ for $n > N$, so sending $\varepsilon \rightarrow 0^+$ as above completes the proof. ■

Proposition 1.75 now gives us quite a bit of control over irrationality measure as long as we have control over the continued fraction. Here are some examples.

Example 1.76. Recall from Example 1.43 that $\varphi = [1; 1, 1, \dots]$. Proposition 1.75 allows us to conclude that $\mu(\alpha) = 0$ immediately because $\log 1 = 0$.

Corollary 1.77. Let $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ be an irrational number with continued fraction $[a_0; a_1, a_2, \dots]$. If there is a polynomial $f \in \mathbb{Z}[x]$ such that $a_n < f(n)$ for sufficiently large n , then $\mu(\alpha) = 2$.

Proof. By Proposition 1.75, it is enough to show that

$$\limsup_{n \rightarrow \infty} \frac{\log a_{n+1}}{\log k_n} = 0.$$

Now, because $a_n < f(n)$ for sufficiently large n , and $f(n) < n^{\deg f + 1}$ for sufficiently large n , it is enough to show

$$\limsup_{n \rightarrow \infty} \frac{d \log n}{\log k_n} \leq 0$$

for any $d > 0$. Of course, we can now factor out the d and thus ignore it.

The main point is that $\{k_n\}_{n=0}^\infty$ increases at least exponentially. Explicitly, we claim is that $k_n \geq 1.5^{n-1}$ for any nonnegative n . This is by induction: certainly $k_0 = 1 \geq 1.5^{-1}$ and $k_1 = a_0 \geq 1.5^0$. Then for the induction, we see

$$k_{n+2} = a_{n+2}k_{n+1} + k_n \geq 1.5^n + 1.5^{n-1} > 1.5^{n-2},$$

where the last inequality holds because it rearranges to $1.5 + 1 > 1.5^2$, which is true.

Applying the main claim, we see that

$$0 \leq \limsup_{n \rightarrow \infty} \frac{\log n}{\log k_n} \leq \limsup_{n \rightarrow \infty} \frac{\log n}{1.5n} = 0,$$

which completes the proof. ■

Corollary 1.78. Let r be any real number at least 2. Then there is an irrational $\alpha \in \mathbb{R} \setminus \mathbb{Q}$ such that $\mu(\alpha) = r$.

Proof. For psychological reasons, we relegate $r = 2$ to Example 1.69. Otherwise, we may assume $r > 2$.

We construct α by infinite continued fraction $[a_0; a_1, a_2, \dots]$, defining the a_n inductively using Proposition 1.32. Define $a_0 := 1$ and $a_1 := 2$ so that we have $k_0 := 1$ and $k_1 := 2$. Then for each $n \geq 1$, define $a_{n+1} := \lfloor k_n^{r-2} \rfloor$ and then $k_{n+1} := a_{n+1}k_n + k_{n-1}$ (as in Proposition 1.32). Then there is an integer sequence $\{h_n\}_{n=0}^\infty$ such that the rational numbers $\{h_n/k_n\}_{n=0}^\infty$ are the continued fraction convergents of $\alpha := [a_0; a_1, a_2, \dots]$. By Corollary 1.48, we see that α is in fact irrational.

It remains to show that $\mu(\alpha) = r$. By Proposition 1.75, we would like to show that

$$r - 2 \stackrel{?}{=} \limsup_{n \rightarrow \infty} \frac{\log a_{n+1}}{\log k_n} = \limsup_{n \rightarrow \infty} \frac{\log \lfloor k_n^{r-2} \rfloor}{\log k_n}.$$

In fact, we claim that

$$\lim_{n \rightarrow \infty} \frac{\log \lfloor n^{r-2} \rfloor}{\log n} \stackrel{?}{=} r - 2,$$

which is good enough upon taking the limit along the subsequence $\{k_n\}_{n=0}^\infty$. Well, we see that

$$r - 2 = \frac{\log n^{r-2}}{\log n} \leq \frac{\log \lfloor n^{r-2} \rfloor}{\log n} \leq \frac{\log n^{r-2}}{\log n} + \frac{1}{\log n} = r - 2 + \frac{1}{\log n}$$

for any sufficiently large n , so we conclude upon taking $n \rightarrow \infty$. ■

1.4.3 Algebraic Bounds on Irrationality Measure

One reason Diophantine approximation attracted the attention of number theorists is that one is able to use the condition that a number is algebraic in order to bound approximations. The prototypical and simplest result of this type is due to Liouville.

Definition 1.79 (algebraic, transcendental). A nonzero complex number $\alpha \in \mathbb{C}$ is *algebraic of degree d* if and only if α is the root of an irreducible polynomial with rational coefficients and of degree d . We denote this degree d by $\deg \alpha$. If no such polynomial exists, then α is called *transcendental*.

Remark 1.80. Let's see that $\deg \alpha$ is well-defined: suppose that α is algebraic and hence the root of an irreducible polynomial f ; by well-ordering, we may choose f to be of least degree. Then for any $g \in \mathbb{Q}[x]$ such that $g(\alpha) = 0$, we claim that f divides g (as polynomials in $\mathbb{Q}[x]$); taking g irreducible then forces $\deg f = \deg g$. Well, using the division algorithm for $\mathbb{Q}[x]$, we may write

$$g(x) = q(x)f(x) + r(x)$$

where $r = 0$ or $0 \leq \deg r < \deg f$. Plugging in α , we see that $r(\alpha) = 0$. Now, if $r \neq 0$, then we may factor r , and one of the irreducible factors will be irreducible, vanish at α , and have degree less than $\deg f$, contradicting the minimality of f . So instead $r = 0$, meaning f divides g .

Proposition 1.81 (Liouville). Fix an algebraic real number $\alpha \in \mathbb{R}$ of degree $d \geq 2$. Then there exists a real number $\varepsilon > 0$ such that

$$\left| \alpha - \frac{h}{k} \right| > \frac{\varepsilon}{k^d}$$

for any rational number h/k with $k > 0$.

Proof. Let f be an irreducible polynomial with integer coefficients of degree $d \geq 2$ where $f(\alpha) = 0$. Namely, we may assume that f has integer coefficients by multiplying out a common denominator.

The main claim is that $|f(h/k)| \geq 1/k^d$. Indeed, $f(h/k) \neq 0$ because f may have no rational roots; explicitly, $f(h/k) = 0$ implies that $kx - h$ divides $f(x)$ by Remark 1.80, but f is irreducible, so this cannot be. But because f has integer coefficients, we may clear denominators to see $k^n f(h/k)$ is an integer. Explicitly, write $f(x) = \sum_{i=0}^d f_i x^i$ for integers f_i , from which we find

$$k^n f\left(\frac{h}{k}\right) = \sum_{i=0}^d f_i h^i k^{d-i} \in \mathbb{Z}.$$

Thus, $|k^n f(h/k)| \geq 1$, and the claim follows.

To see how the above claim helps us, we note that

$$f(\alpha) - f\left(\frac{h}{k}\right) \approx \left(\alpha - \frac{h}{k}\right) f'(\alpha),$$

but the left-hand side has magnitude bounded below by $1/k^d$, so rearranging ought to give the result.

To make this rigorous, we begin by promising $\varepsilon < 1$ so that we might as well assume $|\alpha - h/k| < 1$. This allows us to make \approx above into a genuine equality by using the Mean value theorem to find β between α and h/k such that

$$f(\alpha) - f\left(\frac{h}{k}\right) = \left(\alpha - \frac{h}{k}\right) f'(\beta).$$

Taking absolute values and using the main claim, we find

$$\left| \alpha - \frac{h}{k} \right| \geq \left| \frac{f(h/k)}{f'(\beta)} \right| \geq \frac{1}{|f'(\beta)| k^d}.$$

We now choose $\varepsilon > 0$ small enough so that $|f'(\beta_0)| < 1/\varepsilon$ for any $\beta_0 \in [\alpha - 1, \alpha + 1]$; such an upper bound exists because f' is a continuous function, and $[\alpha - 1, \alpha + 1]$ is compact. Because β is between α and h/k , and h/k is at most 1 away from α , this choice of ε completes the proof. ■

Corollary 1.82. Fix an algebraic real number $\alpha \in \mathbb{R}$ of degree $d \geq 2$. Then $\mu(\alpha) \leq d$.

Proof. For any $\varepsilon > 0$, we show that there are only finitely many rational numbers h/k with $k > 0$ such that

$$\left| \alpha - \frac{h}{k} \right| < \frac{1}{k^{d+\varepsilon}},$$

which will show that $\mu(\alpha) < d+\varepsilon$ and hence complete the proof upon sending $\varepsilon \rightarrow 0^+$. Well, Proposition 1.81 grants some real number $\delta > 0$ such that

$$\left| \alpha - \frac{h}{k} \right| > \frac{\delta}{k^d}$$

for any rational number h/k with $k > 0$. Now, for sufficiently large $k > (1/\delta)^{1/\varepsilon}$, so $1/k^{d+\varepsilon} < \delta/k^d$, so indeed $|\alpha - h/k| < 1/k^{d+\varepsilon}$ is false for sufficiently large k . ■

Example 1.83. By Example 1.72, the number $L := \sum_{n=0}^{\infty} 2^{-n!}$ has $\mu(L) = +\infty$. Thus, Corollary 1.82 implies that L is transcendental.

Historically, examples of the type Example 1.83 were the first numbers proven to be transcendental.

Proposition 1.81 is not at all sharp. Using bivariate polynomials instead of single polynomials, Thue was able to sharpen Proposition 1.81 into the following.

Theorem 1.84 (Thue). Fix an algebraic real number $\alpha \in \mathbb{R}$ of degree $d \geq 3$. Then for any $\varepsilon > 0$, there are only finitely many rational numbers h/k such that

$$\left| \alpha - \frac{h}{k} \right| < \frac{1}{k^{(d+1)/2+\varepsilon}}.$$

In other words, $\mu(\alpha) \leq (d+1)/2$.

The proof of Theorem 1.84 would add between five and ten pages to these notes, so we will not include it. However, it does not require anything much more serious than what we will cover in the remainder of this subsection.

Anyway, Theorem 1.84 is still not sharp. The following result is due to Roth and is work that earned Roth the Fields medal.

Theorem 1.85 (Roth). Fix an algebraic real number $\alpha \in \mathbb{R}$. Then for any $\varepsilon > 0$, there are only finitely many rational numbers h/k such that

$$\left| \alpha - \frac{h}{k} \right| < \frac{1}{k^{2+\varepsilon}}.$$

In other words, $\mu(\alpha) = 2$.

It follows that all the numbers α we constructed in Corollary 1.78 with $\mu(\alpha) > 2$ were in fact transcendental! The proof of Theorem 1.85 would certainly take us too far afield, so we will not show it here.

Even though we will not prove Theorem 1.84, we will use it to the following result on Diophantine equations, as a means to reconnect with our roots.

Theorem 1.86 (Thue). Let $f(x) = \sum_{k=0}^d f_k x^k \in \mathbb{Z}[x]$ be an irreducible polynomial of degree $d \geq 3$. For any $c \in \mathbb{Z}$, the equation

$$\sum_{k=0}^d f_k x^k y^{d-k} = c$$

has finitely many integer solutions (x, y) .

Proof using Theorem 1.84. The idea is that x/y should be a good rational approximation to some real root of f , only finitely many of which should exist by Theorem 1.84.

Suppose for the sake of contradiction that there are infinitely many such solutions $\{(x_n, y_n)\}_{n=0}^\infty$. Our first task is to massage this sequence to converge (in some sense) to a root of f . For any given y , the equation we are solving is a polynomial in x equal to some constant, so there are only finitely many solutions. As such, we must have $|y_n| \rightarrow \infty$ as $n \rightarrow \infty$, so for example we may assume that $y_n \neq 0$ and that $\{|y_n|\}_{n=0}^\infty$ is a strictly increasing sequence. In this case, we see that

$$f\left(\frac{x}{y}\right) = \sum_{k=0}^d f_k \cdot \left(\frac{x}{y}\right)^k = y^{-d} \sum_{k=0}^d f_k x^k y^{d-k} = \frac{c}{y^d}$$

for each solution (x, y) with $y \neq 0$. Now, f is a polynomial of positive degree, so $|f(x)| \rightarrow \infty$ as $x \rightarrow \infty$, so having $|f(x/y)| = c/y^d \leq c$ forces x/y to live in some bounded interval $[-M, M]$. But then the infinite sequence $\{x_n/y_n\}_{n=0}^\infty$ must have a convergent subsequence, so we may assume that $\{x_n/y_n\}_{n=0}^\infty$ does in fact converge. Because $f(x_n/y_n) = c/y_n^d \rightarrow 0$ as $n \rightarrow \infty$, we see that $\{x_n/y_n\}_{n=0}^\infty$ converges to some real root α of f .

Our second task is to bound $|\alpha - x_n/y_n|$. Well, we may factor the irreducible polynomial f over \mathbb{C} as

$$f(x) = f_d \prod_{k=1}^d (x - \alpha_k),$$

where $\{\alpha_1, \dots, \alpha_d\}$ are the roots of f . We go ahead and rearrange the roots so that $\alpha_1 = \alpha$. As usual, note that these roots are disjoint for otherwise any double root would be a root shared by $f(x)$ and $f'(x)$, implying that $\gcd(f'(x), f(x))$ is a nontrivial factor of $f(x)$, thus violating irreducibility. We now see that

$$f_d \prod_{k=1}^d \left| \alpha_k - \frac{x_n}{y_n} \right| = f\left(\frac{x_n}{y_n}\right) = \frac{c}{|y_n|^d}$$

for each $n \geq 2$. We now isolate the $|\alpha - x_n/y_n|$ error term. For each $k \neq 2$, we find

$$\left| \alpha_k - \frac{x}{y} \right| > |\alpha_k - \alpha| - \left| \alpha - \frac{x}{y} \right|.$$

Now, by removing finitely many rational numbers from our sequence $\{x_n/y_n\}_{n=0}^\infty$, we may assume that $|\alpha - x_n/y_n|$ is less than $\frac{1}{2} |\alpha_k - \alpha|$ for each $k \neq 2$, which gives $|\alpha_k - x_n/y_n| > \frac{1}{2} |\alpha_k - \alpha|$. Thus,

$$\left| \alpha - \frac{x_n}{y_n} \right| < \underbrace{\frac{c}{f_d} \prod_{k=2}^d \frac{2}{|\alpha_k - \alpha|}}_{\delta :=} \cdot \frac{1}{|y_n|^d}.$$

Now, for $|y_n|$ sufficiently large, we will have $\delta/|y_n|^d < 1/|y_n|^{(d+1)/2+1/4}$, so the infinitude of these rational approximations $\{x_n/y_n\}_{n=0}^\infty$ is now in direct contradiction with Theorem 1.84. ■

Example 1.87. The polynomial $f(x) := x^3 - 2$ is irreducible of degree 3. So Theorem 1.86 implies that $x^3 - 2y^3 = 10$ has only finitely many integer solutions (x, y) . Indeed, there are at least two integer solutions $(2, -1)$ and $(4, 3)$.

1.4.4 e Is Transcendental

Thus far we have only constructed transcendental numbers by showing they have large irrationality measure, but we know from Proposition 1.73 that almost all real numbers have irrationality measure two. Because the set of algebraic numbers is countable (because the set of integer polynomials is countable), it follows that almost all transcendental numbers have irrationality measure two.

The goal of the next two subsections is to provide a single example of a transcendental number with irrationality measure two. In particular, we will show that e is transcendental and that $\mu(e) = 2$. In this subsection, we will show that e is transcendental; our exposition closely follows [Con]. To give us a flavor of the proof, we begin by showing that e is irrational.

Proposition 1.88. The real number e is irrational.

Proof. The main claim is that there is a sequence of rational numbers $\{p_n/q_n\}_{n=0}^{\infty}$ such that $|q_n e - p_n| \rightarrow 0$ and $q_n \rightarrow \infty$ as $n \rightarrow \infty$. Indeed, we set $q_n := n!$ and $p_n := \lfloor q_n e \rfloor$, which we compute by writing

$$n!e = n! \sum_{k=0}^{\infty} \frac{1}{k!} = \sum_{k=0}^n \frac{n!}{k!} + \sum_{k=n+1}^{\infty} \frac{n!}{k!},$$

so

$$n!e - \sum_{k=0}^n \frac{n!}{k!} = \sum_{k=n+1}^{\infty} \frac{1}{(n+1)(n+2) \cdots (k-1)k} < \sum_{k=n+1}^{\infty} \frac{1}{(n+1)^{k-n}} = \sum_{n=1}^{\infty} \frac{1}{(n+1)^k} = \frac{1/(n+1)}{1 - 1/(n+1)} \leq 1,$$

meaning $p_n = \sum_{k=0}^n n!/k!$, and

$$|q_n e - p_n| < \left| \frac{1/(n+1)}{1 - 1/(n+1)} \right| = \frac{1}{n},$$

so indeed $|q_n e - p_n| \rightarrow 0$ as $n \rightarrow \infty$.

We now complete the proof. Suppose for the sake of contradiction that $e = p/q$ for some rational number p/q with $q > 0$ and $\gcd(p, q) = 1$. Then $|q_n p/q - p_n| \rightarrow 0$ as $n \rightarrow \infty$ by the above argument, so $|q_n p - p_n q| \rightarrow 0$. However, $q_n p - p_n q$ is an integer, so $q_n p = p_n q$ for n sufficiently large. But this does not make sense; for example, choosing n to be any prime, we see that $n \mid q_n$, so $n \mid p_n q$, but $p_n \equiv 1 \pmod{n}$ by definition of p_n , so $n \nmid q$ instead. Thus, q must be larger than any prime, which is a contradiction. ■

Remark 1.89. The above proof is in some sense the same argument as Proposition 1.45 applied to e ; namely, we are using the close approximations p_n/q_n to e in order to derive a contradiction with the fact that all nonzero integers have magnitude at least 1. We have written it in the above manner to make the connection to the following transcendental lemma clearer.

The crux of the above argument is the sequence of rational numbers $\{p_n/q_n\}_{n=0}^{\infty}$ such that $|q_n e - p_n| \rightarrow 0$ as $n \rightarrow \infty$. In order to show that e fails to be algebraic, the key is to find a way to simultaneously approximate not just e but also its powers. The following lemma explains how we will do this approximation.

Lemma 1.90. Fix a nonzero real number $\alpha \in \mathbb{R}$ and a positive integer d . Further, suppose that we have sequences of rational numbers $\{p_{1n}/q_n\}, \{p_{2n}/q_n\}, \dots, \{p_{dn}/q_n\}$ satisfying the following.

- (a) Approximation: for each k , we have $|q_n \alpha^k - p_{kn}| \rightarrow 0$ as $n \rightarrow \infty$.
- (b) Technical: for each n , there is a common divisor g_n of the $p_{\bullet n}$ which is coprime to q_n but satisfies $g_n \rightarrow \infty$ as $n \rightarrow \infty$.

Then α is not the root of an irreducible polynomial in $\mathbb{Z}[x]$ of degree d .

Proof. Suppose for the sake of contradiction that $f(\alpha) = 0$ for some irreducible polynomial $f \in \mathbb{Z}[x]$ of degree d . To be explicit, write $f(x) = a_0 + a_1 x + \cdots + a_d x^d$ where $a_d \neq 0$. Note that $a_0 \neq 0$ because this would require $f(x) = x$, but $\alpha \neq 0$.

Now, the main idea is that p_{kn}/q_n should well-approximate α^k , so we go ahead and plug this into the "linear relation" $f(\alpha) = 0$. For any $n \geq 0$, we write

$$a_0 + \sum_{k=1}^d a_k \cdot \frac{p_{kn}}{q_n} = a_0 + \sum_{k=1}^d a_k \cdot \frac{p_{kn}}{q_n} - f(\alpha) = \sum_{k=1}^d a_k \left(\frac{p_{kn}}{q_n} - \alpha^k \right).$$

Clearing denominators, we find

$$q_n a_0 + \sum_{k=1}^d a_k p_{kn} = - \sum_{k=1}^d a_k (q_n \alpha^k - p_{kn}).$$

As $n \rightarrow \infty$, (a) tells us that the right-hand side goes to 0, so we must have

$$q_n a_0 = - \sum_{k=1}^d a_k p_{kn}$$

for n sufficiently large. However, this cannot be: q_n divides the right-hand side for all n , so $q_n \mid q_n a_0$, so $q_n \mid a_0$, which is a contradiction because a_0 is finite while $q_n \rightarrow \infty$ as $n \rightarrow \infty$. ■

It remains to construct these miraculous rational approximations p_{kn}/q_n of e^k . For this, we must use something about e ; we will input the fact that $\frac{d}{dx} e^x = e^x$ into an integration by parts. To set up the relevant integration by parts, we will define

$$I_f(x) := \sum_{k=0}^{\infty} f^{(k)}(x)$$

for any polynomial f . Notably, this sum is finite because the degree of f is finite. Here is our integration by parts result.

Lemma 1.91. For any polynomial f , we have

$$e^x \int_0^x e^{-t} f(t) dt = e^x I_f(0) - I_f(x).$$

Proof. Quickly, note that $I_f(x)$ is actually a finite sum because f is a polynomial. To get a taste of what is going on, we begin by writing the repeated integration by parts

$$\begin{aligned} \int_0^x e^{-t} f(t) dt &= f(0) - e^{-x} f(x) + \int_0^x e^{-t} f'(t) dt \\ &= (f(0) + f'(0)) - e^{-x} (f(x) + f'(x)) + \int_0^x e^{-t} f''(t) dt. \end{aligned}$$

This process continues. To make this rigorous, we define $I_f^m(x) := \sum_{k=0}^m f^{(k)}(x)$, and we claim that

$$\int_0^x e^{-t} f(t) dt \stackrel{?}{=} I_f^m(0) - e^{-x} I_f^m(x) + \int_0^x e^{-t} f^{(m+1)}(t) dt$$

for any integer $m \geq -1$; the result will follow upon taking $m > \deg f$ so that $f^{(m)} = 0$. We show the claim by induction. At $m = -1$, there is nothing to say. For the inductive step, we note that integration by parts yields

$$I_f^m(0) - e^{-x} I_f^m(x) + \int_0^x e^{-t} f^{(m+1)}(t) dt = I_f^m(0) + f^{(m+1)}(0) - e^{-x} (I_f^m(x) + f^{(m+1)}(x)) + \int_0^x e^{-t} f^{(m+2)}(t) dt,$$

which is what we wanted upon rearranging and plugging into the inductive hypothesis. ■

We are now ready to prove the main result of this subsection.

Theorem 1.92. The real number e is transcendental.

Proof. Note that $e \neq 0$. We will use Lemma 1.90 show that e is not the root of any irreducible polynomial in $\mathbb{Z}[x]$ of degree d , for each $d \geq 1$. Thus, fixing some d , we need to construct the necessary sequences of rational numbers $\{p_{kn}/q_n\}$. For this, we use Lemma 1.91. We would like to approximate e^k , so we plug in $x = k$ to see that

$$e^k \int_0^k e^{-t} f(t) dt = e^k I_f(0) - I_f(k)$$

for any polynomial f . We would like the integral to be relatively small for each k between 0 and d , so we will set

$$f_n(t) := t^{n-1}(t-1)^n(t-2)^n \cdots (t-d)^n$$

for $n \geq 1$. It is also important that f_n vanishes at $k \in \{1, 2, \dots, d\}$ to a higher order than at 0. It is now tempting to directly set p_{kn}/q_n to be $I_{f_n}(k)/I_{f_n}(0)$, but we will want to use the high vanishing of f_n in order to factor out from p_{kn} and q_n beforehand.

Indeed, for each $k \in \{0, 1, \dots, d\}$, we have the Taylor expansion

$$f_n(t+k) = \sum_{\ell=0}^{\infty} \frac{f_n^{(\ell)}(k)}{\ell!} \cdot t^\ell,$$

but these coefficients must all be integers, so we conclude that $\ell! \mid f_n^{(\ell)}(k)$ for all nonnegative integers ℓ . At $k = 0$, we actually have $f_n^{(\ell)}(0) = 0$ for $0 \leq \ell \leq n-1$; and for $k \in \{1, 2, \dots, d\}$, we have $f_n^{(\ell)}(k) = 0$ for $0 \leq \ell \leq n$. Thus, $I_{f_n}(k)$ is divisible by $(n-1)!$ for each k , but it is divisible by $n!$ for each $k > 0$ while

$$I_{f_n}(0) \equiv f^{(n-1)}(0) \equiv (-1)^{nd} d! \pmod{n!} \quad (1.3)$$

because the higher-order terms are 0 (mod $n!$).

With this in mind, we set $p_{kn} := I_{f_n}(k)/(n-1)!$ and $q_n := I_{f_n}(0)/(n-1)!$ for each nonnegative integer n . We now check (a) and (b) of Lemma 1.90, which will complete the proof.

(a) We compute

$$\begin{aligned} |q_n e^k - p_{nk}| &\leq \frac{e^k}{(n-1)!} \int_0^k e^{-t} |f_n(t)| dt \\ &= \frac{e^k}{(n-1)!} \int_0^k e^{-t} |t|^{n-1} |t-1|^n |t-2|^n \cdots |t-d|^n dt \\ &\leq \frac{e^k}{(n-1)!} \cdot d^{n-1+dn} \int_0^k e^{-t} dt. \end{aligned}$$

Now, $\int_0^k e^{-t} dt < \int_0^\infty e^{-t} dt = 1$, so we have the bound

$$|q_n e^k - p_{nk}| < \frac{e^k}{d} \cdot \frac{(d^{d+1})^n}{(n-1)!}.$$

The right-hand side goes to 0 as $n \rightarrow \infty$, so the left-hand side must also.

(b) For this check, we actually want to use a subsequence of the rationals we chose. The common divisor will be $g_n := n$, which we know divides each $p_{kn} = I_{f_n}(k)/(n-1)!$ because $I_{f_n}(k)$ is divisible by $n!$. However, we must verify that there are infinitely many n such that n is relatively prime to q_n . Well, (1.3) implies that it is enough for n to be relatively prime to $d!$, so we may take $n(m) := 1 + md!$ and then take our rationals to be $\{p_{k,n(m)}/q_{n(m)}\}$. This completes the proof. ■

1.4.5 The Continued Fraction of e

In this subsection, we compute the continued fraction expansion of e and then use it to show that $\mu(e) = 2$. Our exposition in this subsection follows [Old70]. We are going to prove that

$$e \stackrel{?}{=} [2; 1, 2, 1, 1, 4, 1, 1, 6, 1, \dots, 1, 2m, 1, \dots].$$

This continued fraction naturally comes in threes, so it will actually be easier to show the related continued fraction

$$\frac{e+1}{e-1} = [2; 6, 10, 14, \dots, 4m+2, \dots].$$

Nonetheless, the main part of our story will unsurprisingly be focused on trying to come up with good rational approximations for e . Anyway, let's jump into a proof.

Proposition 1.93. We have

$$\frac{e+1}{e-1} = [2; 6, 10, 14, \dots, 4m+2, \dots].$$

Proof. For clarity, we proceed in steps.

1. We produce reasonably good rational approximations p_n/q_n (for nonnegative integers n) to e . By Lemma 1.91, we have

$$e \int_0^1 e^{-t} f(t) dt = e I_f(0) - I_f(1)$$

for any polynomial f . We would like to make the integral small in order to produce a good rational approximation of e , so we will take our polynomial to be $f_n(t) := t^n(t-1)^n$. Arguing as in Theorem 1.92, we see that $I_{f_n}(0)$ and $I_{f_n}(1)$ are integers divisible by $n!$. Indeed, the Taylor expansion

$$f_n(t+k) = \sum_{\ell=0}^{\infty} \frac{f_n^{(\ell)}(k)}{\ell!} \cdot t^\ell$$

establishes that $f_n^{(\ell)}(k)$ is an integer divisible by $\ell!$ for any $\ell \geq 0$. However, $f_n^{(\ell)}(k) = 0$ for $k \in \{0, 1\}$ and $0 \leq \ell \leq n$ by construction of f_n , so we conclude that $I_{f_n}(k)$ is divisible by $n!$ for $k \in \{0, 1\}$ because all nonzero terms of the sum

$$I_{f_n}(k) = \sum_{\ell=0}^{\infty} f_n^{(\ell)}(k)$$

are divisible by $n!$.

Thus, we define $q_n := I_{f_n}(0)/n!$ and $p_n := I_{f_n}(1)/n!$. To verify that p_n/q_n is in fact a good rational approximation to e , we write

$$\begin{aligned} |q_n e - p_n| &\leq \frac{e}{n!} \int_0^1 e^{-t} |f_n(t)| dt \\ &= \frac{e}{n!} \int_0^1 e^{-t} |t(t-1)|^n dt \\ &< \frac{e}{n!} \int_0^\infty e^{-t} dt \\ &= \frac{e}{n!}. \end{aligned}$$

2. We produce a recurrence relation for the $\{p_n\}$ and $\{q_n\}$. This will arise purely formally by manipulating the integrals

$$J_{a,b} := \int_0^1 e^{-t} t^a (t-1)^b dt.$$

We have two “moves”: on one hand, integration by parts shows

$$J_{a,b} = \int_0^1 e^{-t} t^a (t-1)^b dt = a \int_0^1 e^{-t} t^{a-1} (t-1)^b dt + b \int_0^1 e^{-t} t^a (t-1)^{b-1} dt = a J_{a-1,b} + b J_{a,b-1} \quad (1.4)$$

for $a, b \geq 1$, and the identity $(t-1)^b = t(t-1)^{b-1} - (t-1)^{b-1}$ shows

$$J_{a,b} = \int_0^1 e^{-t} t^a (t-1)^b dt = \int_0^1 e^{-t} t^{a+1} (t-1)^{b-1} dt - \int_0^1 e^{-t} t^a (t-1)^{b-1} dt = J_{a+1,b-1} - J_{a,b-1} \quad (1.5)$$

for $a \geq 0$ and $b \geq 1$. Now, the main claim of this step is that

$$J_{n,n} \stackrel{?}{=} 2n(2n-1)J_{n-1,n-1} + n(n-1)J_{n-2,n-2} \quad (1.6)$$

for $n \geq 2$. We will prove this using (1.4) and (1.5) repeatedly. Getting started, we write

$$J_{n,n} = nJ_{n-1,n} + nJ_{n,n-1} \quad (1.4)$$

$$= nJ_{n-1,n} + n(J_{n-1,n} + J_{n-1,n-1}) \quad (1.5)$$

$$= 2nJ_{n-1,n} + nJ_{n-1,n-1}.$$

The relation

$$J_{n,n} = 2nJ_{n-1,n} + nJ_{n-1,n-1} \quad (1.7)$$

will be helpful again in a moment. Anyway, we now continue, writing

$$J_{n,n} = nJ_{n-1,n-1} + 2n((n-1)J_{n-2,n} + nJ_{n-1,n-1}) \quad (1.4)$$

$$= (2n^2 + n)J_{n-1,n-1} + (2n^2 - 2n)J_{n-2,n} \\ = (2n^2 + n)J_{n-1,n-1} + (2n^2 - 2n)(J_{n-1,n-1} - J_{n-2,n-1}) \quad (1.5)$$

$$= (4n^2 - n)J_{n-1,n-1} - 2n(n-1)J_{n-2,n-1} \\ = (4n^2 - n)J_{n-1,n-1} - n(J_{n-1,n-1} - (n-1)J_{n-2,n-2}) \quad (1.7)$$

$$= 2n(2n-1)J_{n-1,n-1} + n(n-1)J_{n-2,n-2},$$

which is precisely (1.6).

We now conclude this step. Note that

$$q_n e - p_n = \frac{e}{n!} \int_0^1 e^{-t} t^n (t-1)^n dt = \frac{e}{n!} \cdot J_{n,n},$$

so the recurrence (1.6) implies that

$$q_n e - p_n = 2(2n-1)(q_{n-1} e - p_{n-1}) + (q_{n-2} e - p_{n-2})$$

for $n \geq 2$. Because e is irrational (by Proposition 1.88), collecting terms to put all e s on one side and all integers on the other, we produce the system of recurrences

$$\begin{cases} p_{n+1} = 2(2n+1)p_n + p_{n-1}, \\ q_{n+1} = 2(2n+1)q_n + q_{n-1}, \end{cases} \quad (1.8)$$

for $n \geq 1$. These recurrences essentially explain why the desired continued fraction expansion features $4m+2$.

3. We take a linear combination of the $\{p_n\}$ and $\{q_n\}$ to produce continued fraction convergents. To see why we must do this, we begin by computing (p_0, q_0) and (p_1, q_1) . On one hand, $f_0(t) = 1$, so $I_{f_0}(0) = I_{f_0}(1) = 1$, so $(p_0, q_0) = (1, 1)$. On the other hand, $f_1(t) = t(t-1) = t^2 - t$ had $f'_1(t) = 2t - 1$ and $f''_1(t) = 2$, so $I_{f_1}(0) = 1$ and $I_{f_1}(1) = 3$, so $(p_1, q_1) = (3, 1)$.

The number $p_0/q_0 = 1$ is not a continued fraction convergent of e , so there is no way of shifting our sequences directly in order to produce the continued fraction for e . However, if one sets $(h_{-2}, k_{-2}) = (0, 1)$ and $(h_n, k_n) := ((p_{n+1} + q_{n+1})/2, (p_{n+1} - q_{n+1})/2)$ for each $n \geq -1$, then we see

$$\begin{cases} h_n = 2(2n+1)h_{n-1} + h_{n-2}, \\ k_n = 2(2n+1)k_{n-1} + k_{n-2}, \end{cases}$$

for $n \geq -2$ by an explicit computation at $n = -2$ and (1.8) for $n \geq -1$. Thus, by Proposition 1.32, we see

$$\frac{h_n}{k_n} = [2; 6, 10, 14, \dots, 4n+2]$$

for each $n \geq 0$.

4. We complete the proof. It remains to show that $h_n/k_n \rightarrow (e+1)/(e-1)$ as $n \rightarrow \infty$. The bound

$$|q_n e - p_n| < \frac{e}{n!}$$

now reads

$$|(h_n - k_n)e - (h_n + k_n)| < \frac{e}{(n+1)!},$$

which rearranges to

$$\left| \frac{e+1}{e-1} - \frac{h_n}{k_n} \right| < \frac{e}{(e-1)(n+1)!k_n}.$$

Sending $n \rightarrow \infty$ completes the proof. ■

To produce the continued fraction for e , we need to be able to manipulate continued fractions. We will want the following.

Lemma 1.94. Fix any positive real numbers $a, b, c \in \mathbb{R}$ with $a, b \geq 1$. Then $2/[a; b, c] = 1/[(a-1)/2, 1, 1 + 2/[b-1, a]]$.

Proof. The lower bounds on $a, b, c \in \mathbb{R}$ are merely there to ensure we have no division by zero problems. For example, we have ensured $[b-1; a] > 0$ currently. Anyway, unwrapping, we are trying to show

$$\frac{2}{a + \frac{1}{b + \frac{1}{c}}} \stackrel{?}{=} \frac{1}{\frac{a-1}{2} + \frac{1}{1 + \frac{1}{1 + \frac{2}{b-1 + \frac{1}{c}}}}}.$$

This is purely formal. Taking reciprocals, we are trying to show

$$a + \frac{1}{b + \frac{1}{c}} \stackrel{?}{=} (a-1) + \frac{2}{1 + \frac{1}{1 + \frac{2}{b-1 + \frac{1}{c}}}}.$$

Now, the fraction on the right-hand side is

$$\frac{2}{1 + \frac{1}{1 + \frac{2c}{bc - c + 1}}} = \frac{2}{1 + \frac{bc - c + 1}{bc + c + 1}} = \frac{bc + c + 1}{bc + 1} = 1 + \frac{c}{bc + 1} = 1 + \frac{1}{b + \frac{1}{c}},$$

which completes the proof upon plugging in to the previous equation. ■

Theorem 1.95. We have

$$e = [2; 1, 2, 1, 1, 4, 1, 1, 6, 1, \dots, 1, 2m, 1, \dots].$$

Proof. Subtracting one and taking the reciprocal from Proposition 1.93, we find

$$\frac{e-1}{2} = [0; 1, 6, 10, 14, \dots, 4m+2, \dots].$$

Rearranging, we find

$$e = 1 + 2/[1; 6, 10, 14, \dots, 4m+2, \dots].$$

Beginning our translation, we use Lemma 1.94 to see that this is

$$e = 1 + 1/[0; 1, 1 + 2/[5; 10, 14, 18, \dots]] = [1; 0, 1, 1 + 2/[5; 10, 14, 18, \dots]].$$

More generally, we claim that

$$e \stackrel{?}{=} [1; 0, 1, 1, 2, 1, \dots, 1, 2m, 1, 1 + 2/[4m+5; 4m+10, 4m+14, 4m+18, \dots]]$$

for any $m \geq 0$. We just showed the $m = 0$ case. For the induction, we use Lemma 1.94 to find

$$\begin{aligned} e &= [1; 0, 1, 1, 2, 1, \dots, 1, 2m, 1, 1 + 2/[4m+5; 4m+10, 4m+14, 4m+18, \dots]] \\ &= [1; 0, 1, 1, 2, 1, \dots, 1, 2m, 1, 1 + 1/[2m+1; 1, 1 + 2/[4m+9, 4m+14, 4m+18, \dots]]] \\ &= [1; 0, 1, 1, 2, 1, \dots, 1, 2m, 1, 1, 2m+2, 1, 1 + 2/[4m+9, 4m+14, 4m+18, \dots]]. \end{aligned}$$

Sending $m \rightarrow \infty$ and adjusting the start of the continued fraction completes the proof. Formally, one should justify why sending $m \rightarrow \infty$ makes the continued fraction converge, but this holds essentially by the argument of Proposition 1.40 because the last coefficient of the continued fraction above is always bigger than one and hence unable to cause problems with convergence. ■

Corollary 1.96. We have $\mu(e) = 2$.

Proof. This follows directly from plugging in Theorem 1.95 into Corollary 1.77. For example, the polynomial $f(n) = n + 3$ will do the trick. ■

1.4.6 Problems

Do ten points worth of the following exercises.

Problem 1.4.1 (1 points). Compute the first five continued fraction convergents of e .

Problem 1.4.2 (2 points). Without appeal to results unproven in these notes, work Exercise 1.70.

Problem 1.4.3 (3 points). Show that the real number

$$\sum_{k=0}^{\infty} \frac{n}{10^{n!}}$$

is transcendental.

Problem 1.4.4 (3 points). Compute

$$\int_0^1 e^{-t} t^5 (t-1)^4 dt.$$

Problem 1.4.5 (4 points). Consider the real number

$$L = \sum_{k=0}^{\infty} \frac{1}{2^{3^k}}.$$

Show that $\mu(L) \geq 3$.

Problem 1.4.6 (5 points). For $|y| > 100$, show that any integer pair (x, y) such that $x^3 - 2y^3 = 10$ must have x/y be a continued fraction convergent of $\sqrt[3]{2}$. Using Sage, show that there are no solutions aside $(x, y) \in \{(2, -1), (4, 3)\}$ with $|x|, |y| < 10^{100}$. Please submit the program.

Problem 1.4.7 (10 points). Adapt the proof of Proposition 1.93 to show that

$$\frac{e^{2/k} + 1}{e^{2/k} - 1} = [k; 3k, 5k, \dots]$$

for any integer $k \geq 2$. You may find [Old70] helpful.

THEME 2

QUADRATIC EQUATIONS

2.1 Pell Equations

The goal of the present section is to discuss equations of the form

$$ax^2 + bxy + cy^2 = d$$

where $a, b, c, d \in \mathbb{Z}$. Completing the square and completing the denominator, we might as well solve

$$ax^2 + by^2 = c$$

where $a, b, c \in \mathbb{Z}$. Multiplying through by a , we are trying to solve

$$(ax)^2 + (ab)y^2 = ac,$$

so we may as well try to find integer solutions to the equation

$$x^2 - dy^2 = c$$

where $d, c \in \mathbb{Z}$. If $d < 0$, then $x^2 - dy^2 = c$ must have $|x| \leq \sqrt{c}$ and $|y| < \sqrt{c/d}$, so solving this equation can be done via a finite computation. Otherwise, $d > 0$. From here, it turns out that we can produce much of the internal structure of the solutions by limiting our view to $|c| = 1$, which we will call “Pell equations” for now (though we will want to expand our definition). This will remain our focus for the majority of this section.

2.1.1 Pell Equations via Elementary Methods

Shortly we are going to begin discussing real quadratic fields and their connections to Pell equations, but it is worthwhile to be aware that one can make purely elementary arguments to solve these equations. Let’s see a few examples and feel the wonder.

Remark 2.1. The following examples, suitably transformed, can also be seen as a form of “Vieta jumping.” We will not bother to explain what Vieta jumping is, but those who do may find Problem [2.1.3](#) compelling.

Example 2.2. Define the sequence of ordered pairs of nonnegative integers $\{(x_n, y_n)\}_{n=0}^{\infty}$ recursively by $(x_0, y_0) := (1, 0)$ and

$$(x_{n+1}, y_{n+1}) := (2x_n + 3y_n, x_n + 2y_n)$$

for any $n \geq 0$. Then for any pair of nonnegative integers (x, y) such that $x^2 - 3y^2 = 1$, we have $(x, y) = (x_n, y_n)$ for some nonnegative integer n , and (x_n, y_n) is a solution for each n .

Solution. We have two claims to show, so we will show them separately. The main characters of this solution are the linear transformation $f_{\pm}: \mathbb{Z}^2 \rightarrow \mathbb{Z}^2$ given by $f_{\pm}(x, y) := (2x \pm 3y, \pm x + 2y)$ where \pm is some sign; in particular, $f_{+}(x_n, y_n) = (x_{n+1}, y_{n+1})$.

1. We show that $x_n^2 - 3y_n^2 = 1$ for all nonnegative integers n . We induct on n . At $n = 0$, we are saying $1^2 - 3 \cdot 0^2 = 1$, which is true. For the inductive step, we show that $f_{\pm}(x, y)$ is a solution of (x, y) is for any sign $+$ or $-$. Well, if $x^2 - 3y^2 = 1$, then we compute

$$(2x \pm 3y)^2 - 3(\pm x + 2y)^2 = (4x^2 \pm 12xy + 9y^2) - 3(x^2 \pm 4xy + 4y^2) = x^2 - 3y^2 = 1,$$

so $f(x, y)$ is also a solution.

2. As an intermediate step, we check that f_{+} and f_{-} are inverse functions. Well, we compute

$$f_{\pm}(f_{\mp}(x, y)) = f(2x \mp 3y, \mp x + 2y) = (2(2x \mp 3y) \pm 3(\mp x + 2y), \pm(2x \mp 3y) + 2(\mp x + 2y)) = (x, y)$$

for any arrangement of signs.

3. Lastly, fix a solution (x, y) of $x^2 - 3y^2 = 1$. We would like to show that $(x, y) = (x_n, y_n)$ for some n , which is equivalent to $(x, y) = f_{+}^n(1, 0)$, or $f_{-}^n(x, y) = (1, 0)$ for some n . Well, let n be the largest nonnegative integer such that both entries of $(x', y') := f_{-}^n(x, y)$ are both nonnegative integers; we claim that $(x', y') = (1, 0)$, which will complete the proof. The first step establishes that $(x')^2 - 3(y')^2 = 1$, and we know that one of the coordinates of $f_{-}(x', y')$ must fail to be a nonnegative integer. Thus, we either have $2x' - 3y' < 0$ or $-x' + 2y' < 0$, so $2x' < 3y'$ or $2y' < x'$.

If $2x' < 3y'$, then

$$1 = (x')^2 - 3(y')^2 < \left(\frac{3}{2} \cdot y'\right)^2 - 3(y')^2 < 0,$$

so this cannot be. Thus, we must instead have $x' > 2y'$, from which we find

$$1 = (x')^2 - 3(y')^2 > (2y')^2 - 3(y')^2 = (y')^2,$$

so $y' = 0$, so $(x', y') = (1, 0)$. ■

Example 2.3. Define the sequence of ordered pairs of nonnegative integers $\{(x_n, y_n)\}_{n=0}^{\infty}$ recursively by $(x_n, y_n) := (1, 0)$ and

$$(x_{n+1}, y_{n+1}) := (x_n + 2y_n, x_n + y_n)$$

for any $n \geq 0$. Then for any pair of nonnegative integers (x, y) such that $x^2 - 2y^2 = \pm 1$, we have $(x, y) = (x_n, y_n)$ for some nonnegative integer n , and (x_n, y_n) is a solution for each n . In fact, $x_n^2 - 2y_n^2 = (-1)^n$ for each n .

Solution. Again, we have two claims to show, and we will show them separately. As before, the main characters of this solution are the linear transformations $f_{\pm}: \mathbb{Z}^2 \rightarrow \mathbb{Z}^2$ given by $f_{\pm}(x, y) := (\pm x + 2y, x \pm y)$ where \pm is some sign; in particular, $f_{+}(x_n, y_n) = (x_{n+1}, y_{n+1})$.

1. We show that $x_n^2 - 2y_n^2 = (\pm 1)^n$ for each nonnegative integer n . We proceed by induction. At $n = 0$, we are saying that $1^2 - 2 \cdot 0^2 = 1$. For the inductive step, we suppose $x^2 - 2y^2 = (-1)^n$ and show that $(x', y') := f_{\pm}(x, y)$ has $(x')^2 - 2(y')^2 = -(-1)^{n+1}$. Indeed, we compute

$$(\pm x + 2y)^2 - 2(x \pm y)^2 = (x^2 \pm 4xy + 4y^2) - 2(x^2 \pm 2xy + y^2) = -(x^2 - 2y^2) = (-1)^{n+1}.$$

2. We check that f_{+} and f_{-} are inverse functions. Well, we compute

$$f_{\pm}(f_{\mp}(x, y)) = f_{\pm}(\mp x + 2y, x \mp y) = (\pm(\mp x + 2y) + 2(x \mp y), (\mp x + 2y) \pm (x \mp y)) = (x, y)$$

for any arrangement of signs.

3. Lastly, fix a solution (x, y) of $x^2 - 2y^2 = \pm 1$. We would like to show that $(x, y) = (x_n, y_n)$ for some $n \geq 0$, which is equivalent to $(x, y) = f_+^n(1, 0)$, or $f_-^n(x, y) = (1, 0)$ for some n . Well, let n be the largest nonnegative integer such that both entries of $(x', y') := f_-^n(x, y)$ are both nonnegative integers; we claim that $(x', y') = (1, 0)$, which will complete the proof.

The first step gives $(x')^2 - 2(y')^2 = \pm 1$ still, and because a coordinate of $f_-(x', y')$ must be a negative integer, we have either $2y' < x'$ or $x' < y'$. On one hand, if $x' < y'$, then

$$\pm 1 = (x')^2 - 2(y')^2 < (y')^2 - 2(y')^2 = -(y')^2,$$

so we must have $(y')^2 < \mp 1 \leq 1$, so $y' = 0$, which forces $x' < 0$, which makes no sense. On the other hand, if $2y' < x'$, then

$$\pm 1 = (x')^2 - 2(y')^2 > (2y')^2 - 2(y')^2 = 2(y')^2,$$

so $y' = 0$ is still forced, from which we must have $x' = 1$, so $(x', y') = (1, 0)$. ■

Exercise 2.4. Define the sequence of ordered pairs of nonnegative integers $\{(x_n, y_n)\}_{n=0}^\infty$ recursively by $(x_0, y_0) := (1, 0)$ and

$$(x_{n+1}, y_{n+1}) := (3x_n + 4y_n, 2x_n + 3y_n)$$

for any $n \geq 0$. Then for any pair of nonnegative integers (x, y) such that $x^2 - 2y^2 = 1$, show that $(x, y) = (x_n, y_n)$ for some nonnegative integer n . Describe the solutions to $x^2 - 2y^2 = -1$ similarly.

One might look at Example 2.3 and wonder what all the fuss with $(-1)^n$ is, for it is a perfectly reasonable question to look for solutions to $x^2 - 2y^2 = 1$ on its own, as shown by the above exercise. However, the recursion $(x, y) \mapsto (3x + 4y, 2x + 3y)$ is in some sense “more complicated than it has to be” both in the sense that the coefficients are larger than the recursion in Example 2.3 and also in the sense that this recursion is simply the recursion in Example 2.3 applied twice:

$$(x, y) \mapsto (x + 2y, x + y) \mapsto (3x + 4y, 2x + 3y).$$

As such, it turns out that the “correct” thing to do is in fact to look at solutions to $x^2 - 2y^2 = \pm 1$ and then correct at the end to look at solutions to $x^2 - 2y^2 = 1$. The reason why is explained somewhat but not completely in section 2.1.2. A complete explanation will have to wait for the next section.

The following example provides the extreme end of trying to make our recursions as simple as possible.

Example 2.5. Define the sequence of ordered pairs of nonnegative integers $\{(x_n, y_n)\}_{n=0}^\infty$ recursively by $(x_0, y_0) := (2, 0)$ and

$$(x_{n+1}, y_{n+1}) := \left(\frac{x_n + 5y_n}{2}, \frac{x_{n+1} + y_{n+1}}{2} \right)$$

for any $n \geq 0$. Then for any pair of nonnegative integers (x, y) such that $x^2 - 5y^2 = \pm 4$, we have $(x, y) = (x_n, y_n)$ for some nonnegative integer n , and (x_n, y_n) is a solution for each n . In fact, $x_n^2 + x_n y_n - y_n^2 = (-1)^n \cdot 4$ for each n .

Solution. The proof is similar to the previous two ones. We define the linear transformations $f_\pm : \mathbb{Z}^2 \rightarrow \mathbb{Z}^2$ by $f_\pm(x, y) := \frac{1}{2}(\pm x + 5y, x \pm y)$ so that $f_+(x_n, y_n) = (x_{n+1}, y_{n+1})$.

1. We show that $f_\pm(x, y)$ is a pair of integers of the same parity whenever (x, y) is a pair of integers of the same parity. This verifies that $\{(x_n, y_n)\}_{n=0}^\infty$ is in fact a sequence of integers (by induction) because $(x_0, y_0) = (2, 0)$ is a pair of integers of the same parity. Well, if x and y are both the same parity, then $\pm x + 5y$ and $x \pm y$ are both even, so $f_\pm(x, y)$ is a pair of integers. To see that $\frac{1}{2}(\pm x + 5y)$ and $\frac{1}{2}(x \pm y)$ are both the same parity, we note that

$$\frac{\pm x + 5y}{2} \mp \frac{x \pm y}{2} = 2y \equiv 0 \pmod{2}.$$

2. We show that $x_n^2 - 5y_n^2 = (-1)^n \cdot 4$ for each $n \geq 0$. For $n = 0$, we are saying that $2^2 - 5 \cdot 0^2 = 4$, which is true. For the inductive step, we more generally show that $(x')^2 - 5(y')^2 = \mp 4$ when $(x', y') = f_{\pm}(x, y)$ for $x^2 - 5y^2 = \pm 4$. Well, suppose $x^2 - 5y^2 = \pm 4$, and we compute

$$\left(\frac{\pm x + 5y}{2}\right)^2 - 5\left(\frac{x \pm y}{2}\right)^2 = \frac{x^2 \pm 10xy + 25y^2}{4} - 5 \cdot \frac{x^2 \pm 2xy + y^2}{4} = -(x^2 - 5y^2) = \mp 4.$$

3. We check that f_+ and f_- are inverse functions. Well, we compute

$$\begin{aligned} f_{\pm}(f_{\mp}(x, y)) &= f_{\pm}\left(\frac{\mp x + 5y}{2}, \frac{x \mp y}{2}\right) \\ &= \frac{1}{4}(\pm(\mp x + 5y) + 5(x \mp y), (\mp x + 5y) \pm (x \mp y)) \\ &= (x, y). \end{aligned}$$

4. Lastly, fix a solution (x, y) of $x^2 - 5y^2 = \pm 1$. We would like to show that $(x, y) = f_+^n(2, 0)$ for some nonnegative integer n , which is equivalent to $(2, 0) = f_-^n(x, y)$. As usual, let n be the largest nonnegative integer such that $(x', y') := f_-^n(x, y)$ has coordinates which are nonnegative integers. The second step establishes that $(x')^2 - 5(y')^2 = \pm 4$.

Now, not both coordinates of $f_-(x', y')$ are nonnegative integers, so either $5y' < x'$ or $x' < y'$. On one hand, if $x' < y'$, then

$$1 = (x')^2 - 5(y')^2 < -4(y')^2 \leq 0,$$

which makes no sense. On the other hand, if $5y' < x'$, then

$$1 = (x')^2 - 5(y')^2 > 25(y')^2 - 5(y')^2 = 20(y')^2,$$

so $y' = 0$, so $(x', y') = (2, 0)$, which completes the proof. ■

One can now compute solutions to $x^2 - 5y^2 = \pm 1$ from Example 2.5 by checking which pairs (x_n, y_n) have both coordinates even.

Example 2.6. Define the sequence of ordered pairs of nonnegative integers $\{(x_n, y_n)\}_{n=0}^{\infty}$ recursively by $(x_0, y_0) := (1, 0)$ and

$$(x_{n+1}, y_{n+1}) := (2x_n + 5y_n, x_n + 2y_n)$$

for any $n \geq 0$. Then for any pair of nonnegative integers (x, y) such that $x^2 - 5y^2 = \pm 1$, we have $(x, y) = (x_n, y_n)$ for some nonnegative integer n , and (x_n, y_n) is a solution for each n . In fact, $x_n^2 + x_n y_n - y_n^2 = (-1)^n$ for each n .

Solution. Define (x'_n, y'_n) to be the recursion of Example 2.5; recall that $x'_n \equiv y'_n \pmod{2}$ always. We claim that (x'_n, y'_n) has both terms even if and only if n is divisible by 3. Indeed, applying the recursion three times, we see

$$(x, y) \mapsto \left(\frac{x + 5y}{2}, \frac{x + y}{2}\right) \mapsto \left(\frac{3x + 5y}{2}, \frac{x + 3y}{2}\right) \mapsto (2x + 5y, x + 2y), \quad (2.1)$$

and $x \equiv y \equiv 2x + 5y \equiv x + 2y \pmod{2}$. Thus, the first three terms are $(2, 0)$, $(1, 1)$, and $(3, 1)$, so an induction shows that (x_n, y_n) has both terms even if and only if n is divisible by 3, as needed.

It follows that having $x^2 - 5y^2 = \pm 1$ must have $(x, y) = (x'_{3n}/2, y'_{3n}/2)$ for some nonnegative integer n , and we have $(x'_{3n}/2)^2 - 5(y'_{3n}/2)^2 = (-1)^n$ for each n . To complete the proof, we note that the sequence $\{(x'_{3n}/2, y'_{3n}/2)\}_{n=0}^{\infty}$ has $(x'_0/2, y'_0/2) = (1, 0)$ and recursion given by

$$(x'_{3(n+1)}/2, y'_{3(n+1)}/2) = (2x'_{3n}/2 + 5y'_{3n}/2, x'_{3n}/2 + 2y'_{3n}/2)$$

by the computation in (2.1). The result follows. ■

2.1.2 Pell Equations with Sophistication

Let's explain what's going on in the previous examples. Example 2.2 is quite mysterious because we have removed the context from which

$$f_{\pm}(x, y) = (2x \pm 3y, \pm x + 2y)$$

came from. To explain this, the key is to factor our equation $x^2 - 3y^2 = 1$ into

$$(x - y\sqrt{3})(x + y\sqrt{3}) = 1.$$

Even though we are now working with $\sqrt{3}$ s, this is good because now the problem is completely multiplicative! By a little brute force, one finds the solution $(x, y) = (2, 1)$, so for example $(2 - \sqrt{3})(2 + \sqrt{3}) = 1$. But now the solution $2 + \sqrt{3}$ allows us to find more solutions because $x^2 - 3y^2 = 1$ implies

$$\begin{aligned} 1 &= (x + y\sqrt{3})(x - y\sqrt{3}) \\ &= (x + y\sqrt{3})(2 + \sqrt{3})(x - y\sqrt{3})(2 - \sqrt{3}) \\ &= ((2x + 3y) + (x + 2y)\sqrt{3})((2x + 3y) - (x + 2y)\sqrt{3}). \end{aligned}$$

We now see where f_+ came from, and f_- comes from multiplying out $(x + y\sqrt{3})(2 - \sqrt{3})$ to produce another solution. Note that this also immediately explains why f_+ and f_- are inverse operations with no work: f_+ takes $x + y\sqrt{3}$ and multiplies by $2 + \sqrt{3}$, but then f_- multiplies by $2 - \sqrt{3}$, for a total multiplication by $(2 + \sqrt{3})(2 - \sqrt{3}) = 1$.

One can now translate Example 2.3 into saying that all solutions (x, y) in the nonnegative integers to $x^2 - 3y^2 = 1$ take the form $x + y\sqrt{3} = (2 + \sqrt{3})^n$ for some nonnegative integer n . With this in mind, the argument of Example 2.2 directly generalizes into the following result.

Proposition 2.7. Let d be a non-square positive integer, and suppose that $x^2 - dy^2 = 1$ has a positive integer solution. Let (x_1, y_1) be the minimal positive integer solution in y . Then a pair of integers (x, y) satisfies $x^2 - dy^2 = 1$ if and only if there is a sign $\varepsilon \in \{\pm 1\}$ and an integer $n \in \mathbb{Z}$ such that

$$x + y\sqrt{d} = \varepsilon (x_1 + y_1\sqrt{d})^n.$$

Proof. Before doing anything, we check that all the given pairs (x, y) are in fact solutions. Well, an induction on n shows that

$$x - y\sqrt{d} = \varepsilon (x_1 - y_1\sqrt{d})^n,$$

so

$$x^2 - dy^2 = (x - y\sqrt{d})(x + y\sqrt{d}) = \varepsilon^2 (x_1 - y_1\sqrt{d})^n (x_1 + y_1\sqrt{d})^n = (x_1^2 - dy_1^2)^n = 1.$$

It remains to show that any (x, y) satisfying $x^2 - dy^2 = 1$ have the desired form.

We quickly reduce to the case where $x, y \geq 0$. By adjusting ε , we may assume that $x \geq 0$; it remains to show that $x + y\sqrt{d} = (x_1 + y_1\sqrt{d})^n$ for some $n \in \mathbb{Z}$. Further, if $x^2 - dy^2 = 1$, then $(x - y\sqrt{d})(x + y\sqrt{d}) = 1$, so $(x + y\sqrt{d})^{-1} = x - y\sqrt{d}$. Thus, $x + y\sqrt{d} = (x_1 + y_1\sqrt{d})^n$ if and only if $x - y\sqrt{d} = (x_1 + y_1\sqrt{d})^{-n}$, so we may assume that $y \geq 0$ as well. It remains to show $x + y\sqrt{d} = (x_1 + y_1\sqrt{d})^n$ for some $n \geq 0$.

Because (x_1, y_1) is a solution in positive integers, we see that $x_1 + y_1\sqrt{d} \geq 1 + \sqrt{d} > 1$, so $(x_1 + y_1\sqrt{d})^n$ is an increasing sequence as n varies over nonnegative integers. Thus, for any $x + y\sqrt{d}$, we may find some $n \geq 0$ such that

$$(x_1 + y_1\sqrt{d})^n \leq x + y\sqrt{d} < (x_1 + y_1\sqrt{d})^{n+1},$$

so

$$1 \leq (x + y\sqrt{d}) (x_1 + y_1\sqrt{d})^{-n} < x_1 + y_1\sqrt{d}.$$

We would now like to compare our solutions and use minimality of (x_1, y_1) to conclude. This will require the following result.

Lemma 2.8. Let d be a non-square positive integer. Given nonnegative integer solutions (a_1, b_1) and (a_2, b_2) to $x^2 - dy^2 = 1$, the following are equivalent.

(a) $a_1 + b_1\sqrt{d} \geq a_2 + b_2\sqrt{d}$.

(b) $a_1 \geq a_2$ and $b_1 \geq b_2$.

(c) $a_1 \geq a_2$ or $b_1 \geq b_2$.

The same statements are equivalent once \geq is replaced with $>$.

Proof. Of course, (b) implies (a) and (c). Additionally, if $a^2 - db^2 = 1$, then $a = \sqrt{1 + db^2}$ and $b = \frac{1}{d}\sqrt{a^2 - 1}$, both of which are strictly increasing functions, so (c) implies (b). With this in mind, we note that $a^2 - db^2 = 1$ will imply $a + b\sqrt{d} = \sqrt{1 + db^2} + b\sqrt{d}$, a function strictly increasing in b , so (a) implies (c), completing the proof. ■

Remark 2.9. In light of Lemma 2.8, a solution (x, y) to $x^2 - dy^2 = 1$ in the positive integers which is minimal in y is equivalently minimal in x . (In fact, any reasonable weighting will continue to make (x, y) the smallest solution; see Problem 2.1.2.) As such, we may sloppily say that such a solution is simply “minimal” or “smallest” in the future instead of specifying what it is minimal with respect to.

Now, we define

$$x' + y'\sqrt{d} := (x + y\sqrt{d}) (x_1 - y_1\sqrt{d})^n = (x + y\sqrt{d}) (x_1 + y_1\sqrt{d})^{-n}.$$

We quickly check that $x', y' \geq 0$. The main point is that $x' + y'\sqrt{d} \geq 1$ and $(x')^2 - d(y')^2 = 1$ by construction. For example, $y' = \pm \frac{1}{d}\sqrt{(x')^2 - 1}$, so if $x' < 0$, then $x' + y'\sqrt{d} < x' + \frac{1}{\sqrt{d}}(x') < 0$. Similarly, note that $(x')^2 - d(y')^2 = 1$, so $(x' + y'\sqrt{d}) (x' - y'\sqrt{d}) = 1$, so $y' < 0$ implies that $\sqrt{d} \leq x' - y'\sqrt{d} = 1 / (x' + y'\sqrt{d}) \leq 1$, which does not make sense.

Now, $x', y' \geq 0$ and the inequalities

$$1 \leq x' + y'\sqrt{d} < x_1 + y_1\sqrt{d}$$

enforces $1 \leq x' < x_1$ and $0 \leq y' < y_1$ by Lemma 2.8. Because y_1 is minimal among positive integer solutions, we must have $y' = 0$, so $x' = 1$ is forced. It follows that $x' + y'\sqrt{d} = 1$, so

$$x + y\sqrt{d} = (x_1 + y_1\sqrt{d})^n$$

follows. ■

Remark 2.10. We will show in Proposition 2.13 that the hypothesis that $x^2 - dy^2 = 1$ having a solution is always fulfilled.

Continuing with our examples, Example 2.3 is explained by factoring the equation $x^2 - 2y^2 = \pm 1$ into

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

We can now use the fact that $(1 + \sqrt{2})(1 - \sqrt{2}) = -1$ and thus $(1 + \sqrt{2})(-1 + \sqrt{2}) = 1$ to build the relevant f_{\pm} : we compute

$$(x + y\sqrt{2})(\pm 1 + \sqrt{2}) = (\pm x + 2y, x \pm y),$$

which explains f_+ and f_- , and as a bonus, we still explain why f_+ and f_- are inverse functions.

What made Example 2.3 more complicated than Example 2.2 is that our “smallest solution” $1 + \sqrt{2}$ did not have $(1 + \sqrt{2})(1 - \sqrt{2})$ equal to 1 but instead equal to -1 . We could have worked with $(1 + \sqrt{2})^2 = 3 + 2\sqrt{2}$ instead, but this would in some sense be dishonest: $3 + 2\sqrt{2}$ is not really the smallest solution to an equation of the type $x^2 - dy^2 = c$. One can reasonably ask why

$$x^2 - 3y^2 = -1$$

has no solution, a question answered by looking (mod 4). In contrast, $x^2 - 2y^2 = -1$ has no such local obstruction.

Example 2.5 continues the trend of making the equation more complicated. This time, we want to factor $x^2 - 5y^2 = \pm 4$ as

$$\left(\frac{x + y\sqrt{5}}{2}\right)\left(\frac{x - y\sqrt{5}}{2}\right) = \pm 1,$$

where the point is that our “smallest solution” is given by $\left(\frac{1+\sqrt{5}}{2}\right)\left(\frac{1-\sqrt{5}}{2}\right) = -1$. The presence of these half-integers might seem disconcerting; for example, the product of two numbers of the form $\frac{a}{2} + \frac{b}{2}\sqrt{5}$ need not take that form. However, Example 2.5 finds that we are only interested in numbers of the form $\frac{a}{2} + \frac{b}{2}\sqrt{5}$ where a and b have the same parity, and we can check that

$$\left\{\frac{a}{2} + \frac{b}{2}\sqrt{5} : a, b \in \mathbb{Z} \text{ have the same parity}\right\}$$

is closed under addition and multiplication. As such, we can build f_{\pm} as in Example 2.3: $\left(\frac{1+\sqrt{5}}{2}\right)\left(\frac{-1+\sqrt{5}}{2}\right) = 1$ means we want to compute

$$\left(\frac{x + y\sqrt{5}}{2}\right)\left(\frac{\pm 1 + \sqrt{5}}{2}\right) = \frac{\pm x + 5y}{2} + \frac{x \pm y}{2}\sqrt{5},$$

where $\frac{\pm x + 5y}{2}$ and $\frac{x \pm y}{2}$ are both integers because x and y have the same parity.

There are a number of aspects of the above explanations which are still mysterious, but we will respond to them in time. For example, why did we not consider elements of the form $\frac{a}{2} + \frac{b}{2}\sqrt{3}$ in Example 2.2? For that matter, why not elements of the form $\frac{a}{3} + \frac{b}{3}\sqrt{5}$ in Example 2.5? More fundamentally we keep pulling these small solutions like $2 + \sqrt{3}$ and $1 + \sqrt{2}$ and $\frac{1+\sqrt{5}}{2}$ out from nowhere, so where do they come from? This last question is the one we will focus on answering first.

2.1.3 Using Continued Fractions

In this subsection, we will discuss how to find the smallest solution (x, y) to an equation of the form

$$x^2 - dy^2 = c$$

in some restricted cases. The hope is that one can then use this smallest solution to produce all other solutions as in Examples 2.2, 2.3 and 2.5.

Motivated by the previous section, we factor our equation into

$$(x - y\sqrt{d})(x + y\sqrt{d}) = c.$$

Now, the point is that, if c is small, we need $x - y\sqrt{d}$ to also be abnormally small. Thus, x/y needs to be a good rational approximation of \sqrt{d} . If x/y is good enough, we can use Theorem 1.64 in order to deduce that x/y is a continued fraction convergent of \sqrt{d} . This motivates us to use continued fraction convergents. For example, continued fraction convergents automatically produce small values for $x^2 - dy^2$.

Lemma 2.11. Let $\{h_n/k_n\}_{n=0}^\infty$ be the sequence of continued fraction convergents of \sqrt{d} , where d is a non-square positive integer. Then

$$|h_n^2 - dk_n^2| < 2\sqrt{d} + 1$$

for any $n \geq 0$.

Proof. By Proposition 1.40, we see that

$$\left| \sqrt{d} - \frac{h_n}{k_n} \right| < \frac{1}{k_n^2},$$

so factoring as above yields

$$|h_n^2 - dk_n^2| = |h_n - k_n\sqrt{d}| \cdot |h_n + k_n\sqrt{d}| < \frac{|h_n + k_n\sqrt{d}|}{k_n} = \left| \sqrt{d} + \frac{h_n}{k_n} \right|.$$

To complete our bounding, we use the triangle inequality, writing

$$\left| \sqrt{d} + \frac{h_n}{k_n} \right| \leq 2\sqrt{d} + \left| \sqrt{d} - \frac{h_n}{k_n} \right| < 2\sqrt{d} + \frac{1}{k_n^2} \leq 2\sqrt{d} + 1,$$

so we are done. ■

Remark 2.12. We could alternately recover the above bound by using the bound on s_\bullet given in the proof of Proposition 1.55. The difference here is rather inconsequential.

A pigeonhole argument can use Lemma 2.11 to show that $x^2 - dy^2 = 1$ at least has solutions.

Proposition 2.13. Let d be a positive integer. Then the equation $x^2 - dy^2 = 1$ has a solution in the positive integers.

Proof. Applying the pigeonhole principle to Lemma 2.11, there exists some $N \in \mathbb{Z}$ with $|N| < 2\sqrt{d} + 1$ such that there are infinitely many positive integer solutions (x, y) to $x^2 - dy^2 = N$. Note that there are no positive integer solutions to $x^2 - dy^2 = 0$, so $N \neq 0$. We would like to pick up two solutions (x_1, y_1) and (x_2, y_2) to $x^2 - dy^2 = N$ and write

$$\frac{x_1 + y_1\sqrt{d}}{x_2 + y_2\sqrt{d}}$$

to produce an element with $x^2 - dy^2 = 1$. However, the above element need not take the form $a + b\sqrt{d}$ where a and b are positive integers. Thus, for technical reasons, we note that there are only finitely many elements in $(\mathbb{Z}/|N|\mathbb{Z})^2$, so we must be able to find two distinct pairs of positive integers (x_1, y_1) and (x_2, y_2) such that $x_1 - dy_1^2 = x_2 - dy_2^2 = N$ and $x_1 \equiv x_2 \pmod{N}$ and $y_1 \equiv y_2 \pmod{N}$. (Roughly speaking, this (mod N) business is necessary because of Problem 2.1.1.)

We will now be able to divide $x_1 + y_1\sqrt{d}$ by $x_2 + y_2\sqrt{d}$. To see this, note there exist integers a and b_2 such that

$$x_1 \pm y_1\sqrt{d} = x_2 \pm y_2\sqrt{d} + N(a \pm b\sqrt{d}).$$

Thus,

$$\frac{x_1 \pm y_1\sqrt{d}}{x_2 \pm y_2\sqrt{d}} = 1 + \frac{N}{x_2 \pm y_2\sqrt{d}} \cdot (a \pm b\sqrt{d}) = 1 + (x_2 \mp y_2\sqrt{d})(a \pm b\sqrt{d}).$$

The right-hand side here can thus be expressed as $x \pm y\sqrt{d}$ where x and y are some integers, from which we find

$$(x + y\sqrt{d})(x - y\sqrt{d}) = \left(\frac{x_1 + y_1\sqrt{d}}{x_2 - y_2\sqrt{d}} \right) \left(\frac{x_1 + y_1\sqrt{d}}{x_2 - y_2\sqrt{d}} \right) = \frac{N}{N} = 1.$$

It remains to check that we can coerce (x, y) to live in the positive integers while remaining solutions to $x^2 - dy^2 = 1$. Note that $x = 0$ would imply that $0^2 - dy^2 = 1$, which is impossible. Similarly, $y = 0$ would imply that $x = \pm 1$ and so $(x_2, y_2) = \pm(x_1, y_1)$, which cannot be the case because these are distinct pairs of positive integers. Now, with $x, y \neq 0$, we note that we may adjust their signs to assume that $x, y > 0$ while remaining a solution to $x^2 - dy^2 = 1$, completing the proof. ■

Even though the argument of Proposition 2.13 is somewhat obnoxious, now that we know we have some solution, we can find fairly efficiently using continued fraction convergents, finally executing the approach given at the start of this subsection.

Proposition 2.14. Let d be a non-square positive integer, and suppose (x, y) is a pair of positive integers such that $|x^2 - dy^2| < \sqrt{d}$. Then x/y is a continued fraction convergent of \sqrt{d} .

Proof. We have two cases.

- Suppose $x^2 - dy^2 > 0$. The main point is the bounding

$$\left| \sqrt{d} - \frac{x}{y} \right| = \frac{\left| \frac{x^2}{y^2} - d \right|}{\left| \sqrt{d} + \frac{x}{y} \right|} = \frac{|x^2 - dy^2|}{\left| \sqrt{d} + \frac{x}{y} \right|} \cdot \frac{1}{y^2} < \frac{\sqrt{d}}{\left| \sqrt{d} + \frac{x}{y} \right|} \cdot \frac{1}{y^2}.$$

Now, $x^2 - dy^2 > 0$ implies $x/y > \sqrt{d}$, so we have an upper-bound of $\frac{\sqrt{d}}{2\sqrt{d}} \cdot \frac{1}{y^2} < \frac{1}{2y^2}$, allowing us to conclude by Theorem 1.64.

- Suppose $x^2 - dy^2 < 0$. In this case, we will actually show that y/x is a continued fraction convergent of $1/\sqrt{d}$, which will be enough: if $\sqrt{d} = [a_0; a_1, a_2, \dots]$, then $1/\sqrt{d} = [0; a_0, a_1, a_2, \dots]$, so the reciprocal of a nonzero continued fraction convergent of $1/\sqrt{d}$ will be a continued convergent of \sqrt{d} . Anyway, we bound

$$\left| \frac{1}{\sqrt{d}} - \frac{y}{x} \right| = \left| \frac{x - y\sqrt{d}}{x\sqrt{d}} \right| = \frac{|x^2 - dy^2|}{x\sqrt{d} \cdot |x + y\sqrt{d}|} < \frac{1}{|1 + y/x\sqrt{d}|} \cdot \frac{1}{x^2}.$$

Now, $x^2 - dy^2 < 0$ implies $y/x > 1/\sqrt{d}$, so $|1 + y/x\sqrt{d}| > 2$, so our error is at most $\frac{1}{2x^2}$, allowing us to conclude by Theorem 1.64. ■

Remark 2.15. For example, to find solutions to $x^2 - dy^2 = 1$, we see that we must look among continued fraction convergents $\{h_n/k_n\}_{n=0}^\infty$ of \sqrt{d} , and the periodicity of Remark 1.57 assures us that we might as well check within $0 \leq n \leq 4d$ or so.

Let's use Proposition 2.14 for fun and profit.

Example 2.16. A pair of integers (x, y) satisfies $x^2 - 19y^2 = 1$ if and only if there is a sign $\varepsilon \in \{\pm 1\}$ and an integer $n \in \mathbb{Z}$ such that

$$x + y\sqrt{d} = \varepsilon \left(170 + 39\sqrt{19} \right)^n.$$

Solution. In light of Proposition 2.7, it suffices to show that the smallest solution (x, y) to $x^2 - dy^2 = 1$ is $170 + 39\sqrt{19}$. By Proposition 2.14, it suffices to examine the continued fraction convergents of $\sqrt{19}$ for solutions to $x^2 - dy^2 = 1$. To compute this continued fraction, we use Proposition 1.55. This produces the

(large) table as follows.

n	-2	-1	0	1	2	3	4	5	6	7
r_n			0	4	2	3	3	2	4	4
s_n			1	3	5	2	5	3	1	3
a_n			4	2	1	3	1	2	8	...
h_n	0	1	4	9	13	48	61	170
k_n	1	0	1	2	3	11	14	39

Because $(r_7, s_7) = (r_1, s_1)$, we see that $\sqrt{19} = [4; \overline{2, 1, 3, 1, 2, 8}]$ by the proof of Corollary 1.56. Anyway, we recall from the proof of Proposition 1.55 that $s_n = (-1)^n (h_{n-1}^2 - dk_{n-1}^2)$, so we are looking for $s_n = 1$, for which we see the first nontrivial solution happens at $s_6 = 1$, so $p_5^2 - 19q_5^2 = 1$ is our smallest solution. Referencing the table, this is $(x, y) = (170, 39)$, as needed. ■

Example 2.17. The equation $x^2 - 194y^2 = -1$ has no integer solutions.

Solution. By Proposition 2.14, it suffices to check continued fraction convergents $\{h_n/k_n\}_{n=0}^\infty$ of $\sqrt{194}$, so we use Proposition 1.55 to compute the continued fraction of $\sqrt{194}$, producing the following table.

n	0	1	2	3	4	5
r_n	0	13	12	12	13	13
s_n	1	25	2	25	1	25
a_n	13	1	12	1	26	...

Now, $(r_1, s_1) = (r_5, s_5)$ implies that $(r_{n+4}, s_{n+4}) = (r_n, s_n)$ for any $n \geq 1$ by the recurrence in Proposition 1.55. Because $s_n = (-1)^n (h_{n-1}^2 - 194k_{n-1}^2) > 0$ by the proof of Proposition 1.55 and Corollary 1.56, we see that any solution $h_n^2 - 194k_n^2 = -1$ must have n even while $s_{n+1} = 1$. However, the above periodicity shows that $s_n = 1$ if and only if $n \equiv 0 \pmod{4}$, so n will never be odd. So there are no solutions to $x^2 - 194y^2 = -1$. ■

Remark 2.18. We remark that $13^2 - 194 \cdot 1^2 = -25$ and $5^2 - 194 \cdot 1^2 = -169$, so $(13/5, 1/5)$ and $(5/13, 1/13)$ are both rational solutions to $x^2 - 194y^2 = -1$. We thus claim that $x^2 - 194y^2 \equiv -1 \pmod{n}$ has a solution for each positive integer n , thus breaking a strong form of the local-to-global principle. Indeed, using the Chinese remainder theorem, choose integers (x, y) such that

$$(x, y) \equiv \begin{cases} (13/5, 1/5) \pmod{p^{\nu_p(n)}} & \text{if } p \neq 5, \\ (5/13, 1/13) \pmod{p^{\nu_p(n)}} & \text{if } p = 5. \end{cases}$$

Note that this is in fact finitely many modular conditions because $p^{\nu_p(n)} = 1$ for all primes p except for the p such that $p \mid n$. Anyway, the construction yields $x^2 - dy^2 \equiv -1 \pmod{p^{\nu_p(n)}}$ for all primes p , which yields $x^2 - dy^2 \equiv -1 \pmod{n}$ by the Chinese remainder theorem.

2.1.4 Generalized Pell Equations

Proposition 2.7 combined with the method of Proposition 2.14 tells us how to solve equations of the form

$$x^2 - dy^2 = 1$$

where d is a non-square positive integer. We now use these solutions to solve

$$x^2 - dy^2 = c$$

in general. When $|c| < \sqrt{d}$, we can still use Proposition 2.14, but for larger c , this method does not work. There is still a method to use continued fractions (not of \sqrt{d} but of a similar quadratic irrational) in order to look for solutions, but because we have not discussed an efficient algorithm to compute such continued fractions, we will settle for the following effective result. The proof technique we use below will be used again in various levels of generality.

Proposition 2.19. Fix a non-square positive integer d , and let (x_0, y_0) be a pair of positive integers such that $x_0^2 - dy_0^2 = 1$. Set $u := x_0 + y_0\sqrt{d}$. For any nonzero integer c , any integral solution (x, y) of $x^2 - dy^2 = c$ can be written as $x + y\sqrt{d} = (x' + y'\sqrt{d})u^n$ for some integer n where

$$|x'| \leq \frac{\sqrt{|c|}(\sqrt{u} + 1)}{2} \quad \text{and} \quad |y'| \leq \frac{\sqrt{|c|}(\sqrt{u} + 1)}{2\sqrt{d}}.$$

Proof. Roughly speaking, we would like to measure the height of a nonzero quadratic irrational of the form $x + y\sqrt{d}$. Additionally, we would like for multiplication of the quadratic irrationals to add heights together (to make our arithmetic easier), and we would like for our heights to be positive. It would make sense to set our height, then, to be $\log |x + y\sqrt{d}|$, but from an algebraic point of view, we would like to put \sqrt{d} and $-\sqrt{d}$ on equal footing. As such, we define

$$H(x + y\sqrt{d}) := \left(\log |x + y\sqrt{d}|, \log |x - y\sqrt{d}| \right) \in \mathbb{R}^2,$$

where $x, y \in \mathbb{Q}$ are not both zero. (Note that the value of $x + y\sqrt{d}$ determines x and y because \sqrt{d} is irrational.) Here are some basic properties of H .

- For any $(x, y), (x', y') \in \mathbb{Q}^2 \setminus \{(0, 0)\}$, a direct expansion yields

$$\begin{aligned} (x + y\sqrt{d})(x' + y'\sqrt{d}) &= (xx' + dyy') + (xy' + yx')\sqrt{d} \\ (x - y\sqrt{d})(x' - y'\sqrt{d}) &= (xx' + dyy') - (xy' + yx')\sqrt{d}. \end{aligned}$$

Thus,

$$\begin{aligned} H((x + y\sqrt{d})(x' + y'\sqrt{d})) &= \left(\log |(x + y\sqrt{d})(x' + y'\sqrt{d})|, \log |(x - y\sqrt{d})(x' - y'\sqrt{d})| \right) \\ &= \left(\log |x + y\sqrt{d}|, \log |x - y\sqrt{d}| \right) + \left(\log |x' + y'\sqrt{d}|, \log |x' - y'\sqrt{d}| \right) \\ &= H(x + y\sqrt{d}) + H(x' + y'\sqrt{d}). \end{aligned}$$

For example, an induction implies that $H((x + y\sqrt{d})^n) = nH(x + y\sqrt{d})$ for any $(x, y) \in \mathbb{Q}^2 \setminus \{(0, 0)\}$.

- For any (x, y) such that $x^2 - dy^2 = 1$, we see that $(x + y\sqrt{d})(x - y\sqrt{d}) = 1$, so

$$\log |x + y\sqrt{d}| + \log |x - y\sqrt{d}| = 0.$$

Thus, $H(u^k) \subseteq \{(x, y) \in \mathbb{R}^2 : x + y = 0\}$. In other words, the vectors $H(u)$ and $(1, 1)$ are orthogonal.

More generally, if $a^2 - db^2 = c$ for any nonzero integer c , then $H(a + b\sqrt{d})$ will output to the plane $\{(x, y) \in \mathbb{R}^2 : x + y = \log |c|\}$, which is the set of vectors whose dot product with $(1, 1)$ is $\log |c|$.

Now, suppose that $x^2 - dy^2 = c$, and for brevity let $\alpha := x + y\sqrt{d}$ and $\bar{\alpha} := x - y\sqrt{d}$ so that $\alpha\bar{\alpha} = c$. The point is to estimate the needed n in the proposition by pushing everything through H . Note $(x, y) \neq (0, 0)$ because c is nonzero, so we may place $H(\alpha) \in \mathbb{R}^2$. The second point above tells us that $H(u)$ and $(1, -1)$ form an orthogonal basis of \mathbb{R}^2 , so we write

$$H(\alpha) = sH(u) + t(1, 1),$$

which is

$$(\log |\alpha|, \log |\bar{\alpha}|) = (s \log u + t, -s \log u + t).$$

As roughly explained in the second point above, we can quickly solve for t : summing coordinates, we see that $2t = \log |\alpha \bar{\alpha}|$, so $t = \frac{1}{2} \log |c|$.

We now estimate n as the closest integer to s so that $|n - s| \leq \frac{1}{2}$. This allows us to define $x' + y'\sqrt{d} := \alpha u^{-n}$, and we see that x' and y' are integers because $\alpha = x + y\sqrt{d}$ and $u^{-1} = x_0 - y_0\sqrt{d}$ have all coefficients integers. For brevity, define $\alpha' := x' + y'\sqrt{d}$ and $\bar{\alpha}' := x' - y'\sqrt{d}$. Note $x + y\sqrt{d} = \alpha u^n$ by construction. It remains to prove the inequalities on x' and y' . The idea is to bound α' and $\bar{\alpha}'$ first by passing through H , we see

$$(\log |\alpha'|, \log |\bar{\alpha}'|) = H(\alpha') = H(\alpha) - nH(u) = \left((s - n) \log u + \frac{1}{2} \log |c|, -(s - n) \log u + \frac{1}{2} \log |c| \right).$$

One of $(s - n)$ or $-(s - n)$ is nonnegative. In the case $s - n \geq 0$, then we see that $|\alpha'| = u^{s-n} |c|^{1/2} \leq \sqrt{u |c|}$, and $|\bar{\alpha}'| = u^{-s-n} |c|^{1/2} \leq \sqrt{|c|}$ because $u > 1$. A symmetric pair of inequalities hold in the case where $s - n \leq 0$, so in total we find

$$|\alpha'| + |\bar{\alpha}'| \leq \sqrt{|c|} (\sqrt{u} + 1). \quad (2.2)$$

We are now ready to bound x' and y' . Well, note

$$|x'| = \left| \frac{\alpha' + \bar{\alpha}'}{2} \right| \leq \frac{|\alpha'| + |\bar{\alpha}'|}{2},$$

and

$$|y'| = \left| \frac{\alpha' - \bar{\alpha}'}{2\sqrt{d}} \right| \leq \frac{|\alpha'| + |\bar{\alpha}'|}{2\sqrt{d}},$$

from which (2.2) completes the argument. ■

Example 2.20. A pair of integers (x, y) satisfies $x^2 - 19y^2 = 5$ if and only if there are signs $\varepsilon_x, \varepsilon_y \in \{\pm 1\}$ and an integer $n \in \mathbb{Z}$ such that

$$x + y\sqrt{d} = (\varepsilon_x 9 + \varepsilon_y 2\sqrt{19}) (170 + 39\sqrt{19})^n.$$

Solution. We feed $u := 170 + 39\sqrt{19}$ into Proposition 2.19; recall we found u in Example 2.16. By Proposition 2.19, we would like to find solutions to $x^2 - 19y^2 = 5$ where

$$|y| \leq \frac{\sqrt{5} (\sqrt{170 + 39\sqrt{19}} + 1)}{2\sqrt{19}}.$$

We claim that the right-hand side is less than 6; a more precise calculation could show that it is less than 5, but this will make little difference to us. Squaring it is equivalent to show that

$$\left(\sqrt{170 + 39\sqrt{19}} + 1 \right)^2 < \frac{4 \cdot 19 \cdot 6^2}{5}.$$

Well,

$$\left(\sqrt{170 + 39\sqrt{19}} + 1 \right)^2 < (\sqrt{170 + 40 \cdot 5} + 1)^2 < (\sqrt{400} + 1)^2 = 441 < 456 = 4 \cdot 19 \cdot 6,$$

which is good enough. We now deal with $|y| < 6$ via the following table.

y	$5 + 19y^2$	$x = \pm\sqrt{5 + 19y^2}$
0	5	not integer
1	24	not integer
2	81	± 9
3	176	not integer
4	309	not integer
5	480	not integer

Thus, $(x, y) = (\varepsilon_x 9, \varepsilon_y 2)$ are only solutions with $|y| < 6$. The result follows from Proposition 2.19. ■

Example 2.21. There are no integer solutions (x, y) to $x^2 - 19y^2 = -5$.

Solution. We work as in Example 2.20. Because $|5| = |-5|$, the entire argument goes through until the table. Here is our new table.

y	$-5 + 19y^2$	$x = \pm\sqrt{-5 + 19y^2}$
0	-5	not integer
1	14	not integer
2	71	not integer
3	156	not integer
4	299	not integer
5	470	not integer

Thus, there are no solutions (x, y) to $x^2 - 19y^2 = -5$ in the region $|y| < 6$, which is enough to complete the proof by Proposition 2.19 and the computations of Example 2.20. ■

2.1.5 A Harder Problem

Here is a statement of the same form as Example 2.2.

Proposition 2.22. Define the sequence of ordered triples of nonnegative integers $\{(x_0, y_0, z_0)\}_{n=0}^\infty$ recursively by $(x_0, y_0, z_0) := (1, 0, 0)$ and

$$(x_{n+1}, y_{n+1}, z_{n+1}) = (x_n + 2y_n + 2z_n, x_n + y_n + 2z_n, x_n + y_n + z_n)$$

for any $n \geq 0$. Then for any triple (x, y, z) of nonnegative integers, we have $x^3 + 2y^3 + 4z^3 - 6xyz = 1$ if and only if $(x, y, z) = (x_n, y_n, z_n)$ for some $n \geq 0$.

We do not yet have the machinery to show this result, though it does look like an elementary argument similar to the ones we have already given ought to suffice. A primary goal of our next few weeks is to be able to figure out where statements like Proposition 2.22 come from and provide some idea how to solve them. This will take us through a little algebraic number theory.

2.1.6 Problems

Problem 2.1.1 (2 points). Show that there exist a positive integer d and pairs of positive integers (x_1, y_1) and (x_2, y_2) such that $x_1^2 - dy_1^2 = x_2^2 - dy_2^2$ even though the ratio

$$\frac{x_1 + y_1\sqrt{d}}{x_2 + y_2\sqrt{d}}$$

does not take the form $a + b\sqrt{d}$ where $a, b \in \mathbb{Z}$.

Problem 2.1.2 (2 points). Let d be a non-square positive integer, and fix weights $\alpha, \beta \in \mathbb{R}_{\geq 0}$. Given nonnegative integer solutions (a_1, b_1) and (a_2, b_2) to $x^2 - dy^2 = 1$, show that the following are equivalent.

- (a) $\alpha a_1 + \beta b_1 \geq \alpha a_2 + \beta b_2$.
- (b) $a_1 \geq a_2$.

Problem 2.1.3 (4 points). Define the sequence of ordered pairs of nonnegative integers $\{(x_n, y_n)\}_{n=0}^{\infty}$ recursively by $(x_0, y_0) := (1, 0)$ and

$$(x_{n+1}, y_{n+1}) := (y_n, x_n + y_n)$$

for any $n \geq 0$. Then for any pair of nonnegative integers (x, y) such that $x^2 + xy - y^2 = \pm 1$, show that $(x, y) = (x_n, y_n)$ for some nonnegative integer n .

Problem 2.1.4 (4 points). We study the equation $x^2 - 223y^2 = -3$.

- (a) Show that the equation $x^2 - 223y^2 = -3$ has no integer solutions.
- (b) Show that the equation $x^2 - 223y^2 = -27$ does have integer solutions. Conclude that, for any integer n coprime to 3, the equation $x^2 - 223y^2 \equiv -3 \pmod{n}$ has solutions.

The quickest way to dispatch of the prime 3 is via Problem 3.4.1.

Problem 2.1.5 (5 points). We study the equation $x^2 - 61y^2 = 15$.

- (a) Describe all positive integer solution (x, y) to $x^2 - 61y^2 = 15$.
- (b) Using (a), compute the three smallest positive integers y such that $61y^2 + 15$ is a perfect square. (A little care is required.)

You will almost certainly need to use a computer program; if you use a computer program, please submit it.

Problem 2.1.6 (7 points). Fix a non-square positive integer d , and let $\sqrt{d} = [a_0; a_1, a_2, \dots]$ be the continued fraction with $\{h_n/k_n\}_{n=0}^{\infty}$ as the continued fraction convergents. Let $m \geq 2$ be the least positive integer such that $h_{m-1}^2 - dk_{m-1}^2 = 1$, which exists by Proposition 2.14.

- (a) Using the notation of Proposition 1.55, show that $s_m = 1$ and thus $r_{m+1} = a_0$ and $s_{m+1} = d - a_0^2$.
- (b) Show that $(a_{n+m}, r_{n+m}, s_{n+m}) = (a_n, r_n, s_n)$ for all $n \geq 1$.
- (c) Show that (x, y) is a positive integer solution to $x^2 - dy^2 = 1$ if and only if $(x, y) = (h_{nm-1}, k_{nm-1})$ for some $n \geq 1$.

2.2 Number Rings

In our solutions to Pell equations, we frequently ran into numbers of the form

$$a + b\sqrt{d}$$

where $a, b \in \mathbb{Z}$ and d is a positive integer which is not a square. But in our explanation of Example 2.5 we even ran into numbers of the form

$$\frac{a + b\sqrt{5}}{2}$$

where a and b had the same parity. The goal of this section is to contextualize what is going on here. Starting in this section, we will assume basic ring theory, on the level of any reasonable abstract algebra text. We refer to Appendix A.2 for the necessary field theory.

2.2.1 Normal Domains

It is a question of classical interest in number theory to take an integral domain and ask if it is a unique factorization domain: in some sense, unique factorization domains are “the best” rings (e.g., they are computationally nice to work with because their multiplicative structure can be well-understood). However, it is in general somewhat difficult to do such a check, and one spends a good part of an algebraic number theory learning how to do so.

To start off, a basic hypothesis is that the integral domain be “normal.” For motivation, we recall the Rational root theorem.

Theorem 2.23 (rational root for \mathbb{Z}). Let $f(x) \in \mathbb{Z}[x]$ be a monic polynomial with integer coefficients. If q is a rational root of $f(x)$, then q is an integer.

More generally, we will prove the following more general result.

Theorem 2.24 (rational root). Let A be a unique factorization domain with fraction field K , and let $f(x) \in A[x]$ be a monic polynomial with coefficients in A . If $q \in K$ is a root of $f(x)$, then $q \in A$.

Proof. Quickly, if $q = 0$, there is nothing to say. Otherwise, by using unique factorization, we may write $q = a/b$ where a and b have no irreducible factors in common. Explicitly, by unique factorization, we may express q as a quotient of nonzero elements in A by writing

$$q = \frac{u \prod_{i=1}^n p_i^{\alpha_i}}{\prod_{j=1}^n p_j^{\beta_j}},$$

where u is a unit, α_i and β_i are nonnegative integers, and each p_i is a unique irreducible (not equal to the product of any other p_i and a unit). Then we may write

$$q = u \underbrace{\prod_{\substack{1 \leq i \leq n \\ \alpha_i > \beta_i}} p_i^{\alpha_i - \beta_i}}_{a:=} \bigg/ \underbrace{\prod_{\substack{1 \leq i \leq n \\ \alpha_i < \beta_i}} p_i^{\beta_i - \alpha_i}}_{b:=}$$

Notably, if $\alpha_i = \beta_i$, then we may remove p_i .

We now proceed with the usual proof of the Rational root theorem. Write

$$f(x) = x^d + \sum_{k=0}^{d-1} r_k x^k$$

where $d = \deg f$ and $r_0, r_1, \dots, r_{d-1} \in A$. We are given that $f(a/b) = 0$. As such, the main point is that we can manipulate $f(a/b) = 0$ into

$$0 = a^d + \sum_{k=0}^{d-1} r_k a^k b^{d-k} = a^d + b \sum_{k=0}^{d-1} r_k a^k b^{(d-1)-k}.$$

To show that $q \in A$, we will show that b is a unit, which indeed implies that $q = ab^{-1} \in A$. Because A is a unique factorization domain, it suffices to show that no irreducible element p divides b . Well, if $p \mid b$, then p divides the sum above, so reducing $(\text{mod } p)$ requires $p \mid a^d$, and so $p \mid a$ because p is prime by Lemma A.15; however, no irreducible divides both a and b by their construction, so we are done. ■

Let's see how Theorem 2.24 can detect when a ring fails to be a unique factorization domain.

Example 2.25. Construct the ring $\mathbb{Z}[\sqrt{5}] := \{a + b\sqrt{5} : a, b \in \mathbb{Z}\}$ where addition and multiplication are as expected. Then $\mathbb{Z}[\sqrt{5}]$ is not a unique factorization domain.

Solution. Quickly, we check that $\mathbb{Z}[\sqrt{5}]$ is in fact a subring of (say) \mathbb{C} : we have all the needed identities, and we are closed under the needed operations because

$$\begin{aligned}(a + b\sqrt{5}) + (a' + b'\sqrt{5}) &= (a + a') + (b + b')\sqrt{5}, \\ (a + b\sqrt{5}) \cdot (a' + b'\sqrt{5}) &= (aa' + 5bb') + (ab' + ba')\sqrt{5}.\end{aligned}$$

(We will not check that $\mathbb{Z}[\alpha_1, \dots, \alpha_n]$ is a ring in the future.) Anyway, the main content of this example is to consider the polynomial

$$f(x) := x^2 - x - 1,$$

which is monic and has integer coefficients (in particular, the coefficients are in $\mathbb{Z}[\sqrt{5}]$). Using the quadratic formula, we see that $\frac{1+\sqrt{5}}{2}$ is a root of $f(x)$ which lives in the quotient field of $\mathbb{Z}[\sqrt{5}]$ but not in $\mathbb{Z}[\sqrt{5}]$. Thus, Theorem 2.24 tells us that $\mathbb{Z}[\sqrt{5}]$ cannot be a unique factorization domain! ■

Exercise 2.26. More generally, let $d \equiv 1 \pmod{4}$ be an integer which is not a square. Then show that the set

$$\mathbb{Z}[\sqrt{d}] := \{a + b\sqrt{d} : a, b \in \mathbb{Z}\}$$

is a subring of \mathbb{C} but fails to be a unique factorization domain because $\frac{1+\sqrt{d}}{2}$ is the root of a monic polynomial with integer coefficients.

The test we are applying is worth turning into an adjective.

Definition 2.27 (normal). Fix an integral domain A with fraction field K . Then A is said to be *normal* if and only if the following holds: for any monic polynomial $f(x) \in A[x]$, if $q \in K$ is a root of $f(x)$, then $q \in A$.

Example 2.28. Theorem 2.24 is equivalent to the statement that unique factorization domains are normal.

Non-Example 2.29. The content of Example 2.25 is showing that $\mathbb{Z}[\sqrt{5}]$ is not normal.

It does turn out that the rings $\mathbb{Z}[\sqrt{2}]$ and $\mathbb{Z}[\sqrt{3}]$ are normal, but we will hold off showing this until we have built a little more theory.¹

2.2.2 Number Rings

Notably, we showed that $\mathbb{Z}[\sqrt{5}]$ is not normal by only looking at monic polynomials with coefficients in \mathbb{Z} . This is surprising because the definition of a normal domain allows our coefficients to live in the full ring $\mathbb{Z}[\sqrt{5}]$, but we only used \mathbb{Z} ! It will turn out that using \mathbb{Z} is enough, though showing this requires a little effort. Regardless, we are motivated to make the following definitions.

Definition 2.30 (algebraic integer). Fix a field extension K of \mathbb{Q} . An element $\alpha \in K$ is an *algebraic integer* if and only if α is the root of some monic polynomial in $\mathbb{Z}[x]$.

¹ For example, it is possible to show that these rings are unique factorization domains directly. This approach does not generalize to further $\mathbb{Z}[\sqrt{d}]$ because, for example, $\mathbb{Z}[\sqrt{15}]$ is normal but not a unique factorization domain.

Definition 2.31 (number ring). Fix a finite field extension K of \mathbb{Q} . Then the *number ring* \mathcal{O}_K of K consists of all the algebraic integers in K .

Example 2.32. We see that $\mathcal{O}_{\mathbb{Q}} = \mathbb{Z}$ by Theorem 2.24. Explicitly, any element $n \in \mathbb{Z}$ is the root of the monic polynomial $x - n \in \mathbb{Z}[x]$. On the other hand, if $\alpha \in \mathbb{Q}$ is an algebraic integer, then α is the root of a monic polynomial in $\mathbb{Z}[x]$, so $\alpha \in \mathbb{Z}$ by Theorem 2.24.

Though we will not dwell on it too much, it is worth acknowledging that the following generalized (relative) notion is more correct to work with.

Definition 2.33 (integral). Fix an embedding of rings $A \subseteq B$. An element $\alpha \in B$ is *integral over A* if and only if α is the root of some monic polynomial in $A[x]$.

Example 2.34. Fix a field extension K of \mathbb{Q} . Then $\alpha \in K$ is an algebraic integer if and only if α is integral over \mathbb{Z} .

We do need to check that \mathcal{O}_K is in fact a ring. Of course, 0 and 1 are algebraic integers (they are the roots of the monic polynomials x and $x - 1$, respectively), so it remains to show that \mathcal{O}_K is closed under addition and multiplication. This requires a little work. The following result is known as the “determinant trick” in commutative algebra. Note that this result (and its corollaries) is basically Proposition A.29.

Proposition 2.35. Fix an embedding of rings $A \subseteq B$ and some $\alpha \in B$. Then the following are equivalent.

- (a) α is integral over A .
- (b) The ring $A[\alpha]$ is finitely generated as an A -module.
- (c) There is a subring $A' \subseteq B$ finitely generated as an A -module containing both A and α .

Here, a ring A' containing A is finitely generated as an A -module (or “finitely generated over A ”) if and only if there are finitely many elements r'_1, \dots, r'_n such that each $r' \in A'$ can be written as

$$r' = \sum_{k=1}^n a_k r'_k$$

for some $a_1, \dots, a_n \in A$. In other words, each element of A' is an A -linear combination of some fixed finite set in A' .

Proof. We show the implications separately.

- (a) To see that (a) implies (b), we suppose α is the root of the monic polynomial $f(x) \in A[x]$ of degree d , written as

$$f(x) = x^d + \sum_{k=0}^{d-1} a_k x^k.$$

Then, for any $n \geq d$, we can express α^n as a \mathbb{Z} -linear combination of lower powers because

$$\alpha^n = - \sum_{k=0}^{d-1} a_k \alpha^{k+d-n}.$$

It follows that $A[\alpha]$ is generated by the elements $1, \alpha, \alpha^2, \dots, \alpha^{d-1}$.

- (b) Note that (b) implies (c) by setting $A' = A[\alpha]$.

- (c) Checking that (c) implies (a) is harder. If $A' = 0$, then $\alpha = 0$, so there is nothing to say; otherwise, $A' \neq 0$. Suppose A' is generated by the elements r'_1, r'_2, \dots, r'_n . Note that $\alpha r'_i \in A'$ for each r'_i , so we may write

$$\alpha r'_i = \sum_{j=1}^n a_{ij} r'_j$$

for some elements $a_{ij} \in A$. In other words, the matrix $T := (a_{ij})_{i,j=1}^n$ has

$$\alpha \begin{bmatrix} r'_1 \\ \vdots \\ r'_n \end{bmatrix} = T \begin{bmatrix} r'_1 \\ \vdots \\ r'_n \end{bmatrix}.$$

Thus, $T - \alpha I_n$ is an $n \times n$ matrix with entries in A' , and it has the nonzero vector (r'_1, \dots, r'_n) in its kernel, so $\det(T - \alpha I_n) = 0$. Expanding out the polynomial $\det(\alpha I_n - T) = 0$ makes α the root of a monic polynomial (of degree n) with coefficients in A , so α is indeed integral over A . ■

In our application, we will also want the following lemma.

Lemma 2.36. Let $A \subseteq B \subseteq C$ be embeddings of rings. If C is finitely generated as a B -module, and B is finitely generated as an A -module, then C is finitely generated as an A -module.

Proof. The point is to concatenate our generating sets together; indeed, this argument is basically the same as Lemma A.20. Say that B is generated over A by the elements b_1, \dots, b_m , and say that C is generated over B by the elements c_1, \dots, c_n . Then any $c \in C$ has some $b'_1, \dots, b'_n \in B$ such that

$$c = \sum_{\ell=1}^n b'_\ell c_\ell,$$

but now each b'_ℓ can be expanded over A as

$$c = \sum_{\ell=1}^n \sum_{k=1}^m a_{k\ell} b_k c_\ell$$

for some $a_{k\ell} \in A$. Thus, the elements $b_k c_\ell$ for $1 \leq k \leq m$ and $1 \leq \ell \leq n$ generate C over A , so we are done. ■

Corollary 2.37. Fix a finite field extension K of \mathbb{Q} . Then \mathcal{O}_K is a normal ring.

Proof. We run our checks separately.

- We check that \mathcal{O}_K is a ring. As discussed previously, $0, 1 \in \mathcal{O}_K$ because these elements are the roots of the polynomials x and $x - 1$, respectively. It remains to show that, for any $\alpha, \beta \in \mathcal{O}_K$, we have $\alpha + \beta, \alpha\beta \in \mathcal{O}_K$. The main point is to show that $\mathbb{Z}[\alpha, \beta]$ is finitely generated as a \mathbb{Z} -module, which will complete the proof by Proposition 2.35 because $\alpha + \beta, \alpha\beta \in \mathbb{Z}[\alpha, \beta]$.

Well, let α and β be the roots of the monic polynomials $f(x), g(x) \in \mathbb{Z}[x]$ respectively. Then by Proposition 2.35 shows that $\mathbb{Z}[\beta]$ is finitely generated as a \mathbb{Z} -module, and $f(\alpha) = 0$ shows that α is integral over $\mathbb{Z}[\beta]$, so $\mathbb{Z}[\alpha, \beta]$ is finitely generated as a $\mathbb{Z}[\beta]$ -module. We conclude $\mathbb{Z}[\alpha, \beta]$ is finitely generated as a \mathbb{Z} -module by Lemma 2.36.

- We check that \mathcal{O}_K is normal. Suppose that $\alpha \in K$ is the root of the monic polynomial $f(x) \in \mathcal{O}_K[x]$; we show that $\alpha \in \mathcal{O}_K$ by showing that α is an algebraic integer. Well, expand $f(x)$ as

$$f(x) = x^d + \sum_{k=0}^{d-1} a_k x^k$$

for some $a_0, \dots, a_{d-1} \in \mathcal{O}_K$. Each a_\bullet is integral over \mathbb{Z} , so Proposition 2.35 tells us that $\mathbb{Z}[a_\bullet]$ for each a_\bullet . As such, as in the previous check, we may build the tower

$$\mathbb{Z} \subseteq \mathbb{Z}[a_0] \subseteq \mathbb{Z}[a_0, a_1] \subseteq \dots \subseteq \mathbb{Z}[a_0, \dots, a_{d-1}],$$

where each ring is finitely generated over the previous one by Proposition 2.35. Then Lemma 2.36 tells us that $\mathbb{Z}[a_0, \dots, a_{d-1}]$ is finitely generated as a \mathbb{Z} -module. Lastly, $f(\alpha) = 0$ tells us that α is integral over $\mathbb{Z}[a_0, \dots, a_{d-1}]$, so $\mathbb{Z}[a_0, \dots, a_{d-1}, \alpha]$ is finitely generated as a $\mathbb{Z}[a_0, \dots, a_{d-1}]$ -module and hence finitely generated as a \mathbb{Z} -module by Lemma 2.36, meaning that α is integral over \mathbb{Z} by Proposition 2.35. ■

After doing all that theory, we are owed an example, so we will compute \mathcal{O}_K for quadratic field extensions $K = \mathbb{Q}(\sqrt{d})$ of \mathbb{Q} . We will want the following lemma, which we have been using in some guise for quite a bit of the previous section.

Lemma 2.38. Fix $K := \mathbb{Q}(\sqrt{d})$ where d is a non-square integer. Then the function $\sigma: K \rightarrow K$ given by $\sigma(a + b\sqrt{d}) := a - b\sqrt{d}$ is a ring homomorphism.

Proof. To begin, note that σ is well-defined because $a + b\sqrt{d} = a' + b'\sqrt{d}$ implies that $(a - a') = (b' - b)\sqrt{d}$ and so $a = a'$ and $b = b'$ because \sqrt{d} is irrational. To check that σ is a homomorphism, we note that $\sigma(0) = 0$ and $\sigma(1) = 1$ and

$$\begin{aligned} \sigma((a + b\sqrt{d}) + (a' + b'\sqrt{d})) &= (a + a') - (b + b')\sqrt{d} \\ &= \sigma(a + b\sqrt{d}) + \sigma(a' + b'\sqrt{d}) \\ \sigma((a + b\sqrt{d})(a' + b'\sqrt{d})) &= (aa' + dbb') - (ab' + ba')\sqrt{d} \\ &= \sigma(a + b\sqrt{d})\sigma(a' + b'\sqrt{d}), \end{aligned}$$

completing the proof. ■

Example 2.39. Fix $K := \mathbb{Q}(\sqrt{2})$. Then $\mathbb{Z}[\sqrt{2}] = \mathcal{O}_K$. In particular, $\mathbb{Z}[\sqrt{2}]$ is normal.

Proof. Note that $\sqrt{2} \in \mathcal{O}_K$ because it is the root of the polynomial $x^2 - 2 = 0$. Thus, $\mathbb{Z}[\sqrt{2}]$ is finitely generated as a \mathbb{Z} -module by Proposition 2.35, and we see $\mathbb{Z}[\sqrt{2}] \subseteq \mathcal{O}_K$.

We now show that $\mathcal{O}_K \subseteq \mathbb{Z}[\sqrt{2}]$, which is harder. Suppose $a + b\sqrt{2} \in \mathcal{O}_K$, where we allow $a, b \in \mathbb{Q}$. We want to show that $a, b \in \mathbb{Z}$. Well, $a + b\sqrt{2}$ is the root of some monic polynomial $f(x) \in \mathbb{Z}[x]$, so $\sigma(a + b\sqrt{2}) = a - b\sqrt{2}$ is also the root of $f(x)$ by Lemma 2.38, so $a - b\sqrt{2}$ is also an algebraic integer. Thus,

$$\begin{aligned} (a + b\sqrt{2}) + (a - b\sqrt{2}) &= 2a, \\ ((a + b\sqrt{2}) - (a - b\sqrt{2}))\sqrt{2} &= 4b, \\ (a + b\sqrt{2})(a - b\sqrt{2}) &= a^2 - 2b^2, \end{aligned}$$

are also algebraic integers. But they are also rational and hence actually integers by Example 2.32, so we may write $a = a_0/2$ and $b = b_0/4$ for some integers a_0 and b_0 . Then

$$a^2 - 2b^2 = \frac{a_0^2}{4} - \frac{b_0^2}{8} = \frac{2a_0^2 - b_0^2}{8}$$

needs to be an integer, so $2a_0^2 \equiv b_0^2 \pmod{8}$, which can only happen if $a_0 \equiv 0 \pmod{2}$ and $b_0 \equiv 0 \pmod{4}$. Thus, a and b are in fact integers, so we conclude. ■

More generally, we can show the following statement.

Proposition 2.40. Fix a non-square squarefree integer d , and fix $K := \mathbb{Q}(\sqrt{d})$. Then we have

$$\mathcal{O}_K = \begin{cases} \mathbb{Z}[\sqrt{d}] & \text{if } d \equiv 2, 3 \pmod{4}, \\ \mathbb{Z}\left[\frac{1+\sqrt{d}}{2}\right] & \text{if } d \equiv 1 \pmod{4}. \end{cases}$$

Proof. We proceed as in the example. Of course, $\sqrt{d} \in \mathcal{O}_K$ because it is the root of $x^2 - d = 0$, and if $d \equiv 1 \pmod{4}$, then $\frac{1+\sqrt{d}}{2} \in \mathcal{O}_K$ because it is the root of

$$x^2 - x - \frac{d-1}{4} = 0.$$

Thus, Proposition 2.35 then assures us that $\mathbb{Z}[\sqrt{d}] \subseteq \mathcal{O}_K$ and $\mathbb{Z}\left[\frac{1+\sqrt{d}}{2}\right] \subseteq \mathcal{O}_K$ when $d \equiv 1 \pmod{4}$.

It remains to show our other inclusion. Well, fix some $a + b\sqrt{d} \in \mathcal{O}_K$ where $a, b \in \mathbb{Q}$. Because $a + b\sqrt{d}$ is the root of some monic polynomial with integer coefficients, we see that $a - b\sqrt{d}$ is as well by Lemma 2.38, so $a - b\sqrt{d}$. Thus,

$$\begin{aligned} (a + b\sqrt{d}) + (a - b\sqrt{d}) &= 2a, \\ \left((a + b\sqrt{d}) - (a - b\sqrt{d})\right) \sqrt{d} &= 2bd, \\ (a + b\sqrt{d})(a - b\sqrt{d}) &= a^2 - db^2, \end{aligned}$$

are also algebraic integers. But they are also rational and hence actually integers by Example 2.32, so we may write $a = a_0/2$ and $b = b_0/(2d)$ for some integers a_0 and b_0 . For example, we see that

$$4(a^2 - db^2) = a_0^2 - \frac{b_0^2}{d}$$

must be an integer, so because d is squarefree, we conclude that $d \mid b_0$, so we write $b = b_1/2$ for some integer b_1 . To continue the argument, we split into cases.

- Suppose $d \equiv 2, 3 \pmod{4}$. Then we see

$$a^2 - db^2 = \frac{a_0^2 - db_1^2}{4}$$

must be an integer. By checking $\pmod{4}$, we see that both a_0 and b_1 must be even, so a and b are both integers, so $a + b\sqrt{d} \in \mathbb{Z}[\sqrt{d}]$.

- Suppose $d \equiv 1 \pmod{4}$. Then we see

$$a^2 - db^2 = \frac{a_0^2 - db_1^2}{4}$$

must be an integer. By checking $\pmod{4}$, we see that both a_0 and b_0 must have the same parity, so

$$a + b\sqrt{d} = \frac{a_0 - b_0}{2} + a_0 \cdot \frac{1 + \sqrt{d}}{2} \in \mathbb{Z}\left[\frac{1 + \sqrt{d}}{2}\right],$$

establishing the needed inclusion. ■

Note that Proposition 2.40 fulfills a curiosity of Example 2.5, namely about where the denominator of 2 came from and why it was so controlled!

To continue our journey of generalization to compute other rings of integers, we need to generalize aspects of the above proof: where did the σ come from? Where did the various $2a$, $2bd$, and $a^2 - db^2$ come from? Why were we able to come up with a denominator of $2d$ so quickly, and why was it so annoying to argue beyond that?

2.2.3 The Discriminant**2.2.4 Number Ring Structure****2.2.5 Dirichlet's Unit Theorem: Upper Bound****2.3 Minkowski Theory**

Having spent a long time building the theory of number rings, we will take a break to discuss some geometry of numbers. We will then return to complete the proof of Dirichlet's unit theorem.

2.3.1 Minkowski's Theorem**2.3.2 Dirichlet's Unit Theorem: Upper Bound****2.4 Binary Quadratic Forms**

INTERMISSION: OTHER FIELDS

3.1 Cyclotomic Extensions

3.2 (Almost) Unique Factorization

3.3 Local Fields

3.4 Hensel's Lemma

Problem 3.4.1. Show that $x^2 - 223y^2 = -3$ has a solution in \mathbb{Z}_3 by showing that there exists $y \in \mathbb{Z}_3$ such that $y^2 = 4/223$. Use Problem 2.1.4 to conclude that the equation

$$x^2 - 223y^2 \equiv -3 \pmod{n}$$

has a solution for any positive integer n .

THEME 4

CUBIC EQUATIONS

Every person believes that he knows what a curve is until he has learned so much mathematics that the countless possible abnormalities confuse him.

—Felix Klein, [Kle16]

- 4.1 Elliptic Curves**
- 4.2 Torsion of Elliptic Curves**
- 4.3 Elliptic Curves over Finite Fields**
- 4.4 Modern Perspectives**

APPENDIX A

SOME ALGEBRA

A.1 Unique Factorization Domains

The goal of the present section is to review the notion of a unique factorization domain and basic properties of them. This is a notion we will want later in section 3.2, though it is not clear why yet.

Definition A.1 (integral domain). A ring A is an *integral domain* if and only if $a \cdot b = 0$ implies that $a = 0$ or $b = 0$.

Example A.2. The ring \mathbb{Z} is an integral domain. Any field is an integral domain.

The best integral domains are unique factorization domains. To recall this definition, we need the notion of prime and irreducible elements.

Definition A.3 (prime). Fix a ring A . Then an element $p \in A$ is *prime* if and only if p is nonzero and the ideal (p) of A is a prime ideal. In other words, p is not zero, not a unit, and whenever $p \mid ab$ for $a, b \in A$, we have $p \mid a$ or $p \mid b$.

Definition A.4 (irreducible). Fix a ring A . Then an element $p \in A$ is *irreducible* if and only if any factorization $p = ab$ has exactly one of a or b equal to a unit. Notably, p cannot be zero (for we could set $a = b = 0$), and p cannot be a unit (for then both a and b would be a unit).

Definition A.5 (unique factorization domain). Fix an integral domain A . Then A is a *unique factorization domain* if and only if any nonzero element $r \in A$ has a unique factorization into irreducibles

$$r = \prod_{i=1}^n p_i,$$

where the sequence $\{p_i\}_{i=1}^n$ of irreducible elements of A is unique up multiplication by a unit and permutation.

An elementary number theory course would show that an integer is prime if and only if it is irreducible and then deduce that \mathbb{Z} is a unique factorization domain. An algebra course would show the more general result that a principal ideal domain is a unique factorization domain. We will quickly review these arguments, but we will not dwell on them.

Proposition A.6. The ring \mathbb{Z} is a principal ideal domain.

Proof. Let $I \subseteq \mathbb{Z}$ be an ideal. If $I = \{0\}$, then $I = (0)$. Otherwise, I contains a nonzero element $n \in I$, so I contains a positive element $n^2 \in I$. Thus, we may let $g \in I$ denote the least positive element. We claim that $I = (g)$; certainly $(g) \subseteq I$, so we want to show $I \subseteq (g)$. Well, choose any $a \in I$. Then we may use division to find integers $q, r \in \mathbb{Z}$ such that

$$a = gq + r$$

where $0 \leq r < g$. Thus, $r = a - gq \in I$ is a nonnegative element of I strictly less than g , so r cannot be positive, so $r = 0$, so $a = gq$, so $a \in (g)$. ■

We now move towards showing that principal ideal domains are unique factorization domains.

Lemma A.7. Fix an integral domain A . If p is prime, then p is irreducible.

Proof. Note that p is neither zero nor a unit by hypothesis. Suppose we factor $p = ab$ where $a, b \in A$; certainly both a and b cannot both be units because then p would be a unit, so it remains to show that one is. Then $p \mid ab$, so $p \mid a$ or $p \mid b$ because p is prime. Without loss of generality, take $p \mid a$, and write $a = pa'$ so that

$$p = pa'b.$$

Then $1 = a'b$, so b is a unit. ■

Lemma A.8. Let A be a principal ideal domain. If $p \in A$ is irreducible, and if $a \in A$ lives outside (p) , then $(a, p) = A$. In other words, (p) is a maximal ideal.

Proof. Note that the last sentence follows from the previous because (p) would then be proper ideal with no ideal between (p) and A , making (p) maximal.

Now, note (a, p) is a principal ideal, so say $(a, p) = (d)$. However, because d divides p , we may write $p = de$ where e is an integer. Thus, one of d or e is a unit. We claim that d is a unit, which will complete the proof. Well, if e is a unit, then $d = pe^{-1}$, so pe^{-1} divides a , so p divides a (recall e is a unit), which is a contradiction. ■

Remark A.9. One can interpret Lemma A.8 as showing that $A/(p)$ is a field for any irreducible p . Indeed, this follows from (p) being a maximal ideal.

Proposition A.10. Let A be a principal ideal domain. Then any irreducible element $p \in A$ is prime.

Proof. Note that p is neither zero nor a unit by hypothesis. Now, suppose we have $p \mid ab$ but $p \nmid a$. We want to show that $p \mid b$. Well, Lemma A.8 tells us that $(a, p) = A$, so we can write

$$ar + ps = 1$$

for some integers $r, s \in A$, so we see that

$$abr + psb = b,$$

so $p \mid ab$ and $p \mid psb$ implies that $p \mid b$, which is what we wanted. ■

Theorem A.11. Any principal ideal domain A is a unique factorization domain.

We will split the proof of Theorem A.11 into a few parts. To begin, we prove existence, which does not require Proposition A.10.

Lemma A.12. Fix an integral domain A . Suppose that any ascending chain of principal ideals

$$(a_0) \subseteq (a_1) \subseteq (a_2) \subseteq \cdots$$

eventually stabilizes; in other words, there is some nonnegative integer N such that $(a_n) = (a_N)$ for any $n \geq N$. Then every nonzero element of A has a factorization into irreducibles.

Proof. Units take the “empty” factorization consisting of only the unit itself and no irreducibles. Irreducible elements attain the factorization consisting of the irreducible itself.

Now, suppose for the sake of contradiction that we have a nonzero element $r_0 \in A$ with no factorization into irreducibles. The work above shows that r_0 is not a unit and not irreducible, so it follows that we can factor $r_0 = s_1 r_1$ where neither s_1 nor r_1 is a unit or zero. If s_1 and r_1 both had factorizations into irreducibles, then we could multiply the factorizations together to produce a factorization for r_0 . Thus, at least one of s_1 or r_1 cannot have a factorization into irreducibles; without loss of generality, it is r_1 .

Iterating the process of the previous paragraph produces a sequence of elements $\{r_n\}_{n=0}^{\infty}$ and $\{s_n\}_{n=1}^{\infty}$ such that $r_n = r_{n+1} s_{n+1}$ for each n . But then we have the descending chain of principal ideals

$$(r_0) \subseteq (r_1) \subseteq (r_2) \subseteq \cdots,$$

which must stabilize eventually. Thus, there is some n for which $(r_n) = (r_{n+1})$, so we may find s' such that $r_{n+1} = s' r_n$ and thus

$$r_n = s_n r_{n+1} = s_n s' r_n.$$

This implies that $s_n s' = 1$ and so s_n is a unit, which is a contradiction to its construction. ■

Remark A.13. More generally, a ring A will be called “Noetherian” if and only if any ascending chain of ideals stabilizes. We will avoid using this notion in these notes when possible, largely because it is not strictly necessary for the story we wish to tell.

Lemma A.14. Fix an integral domain A . Suppose that an element $p \in A$ is prime if and only if it is irreducible. Then for any equal factorizations of irreducibles

$$\prod_{i=1}^m p_i = \prod_{j=1}^n q_j,$$

we must have $m = n$, and there is a permutation σ of $\{1, 2, \dots, n\}$ such that p_i and $q_{\sigma(i)}$ are the same up to multiplication by a unit.

Proof. Fix a factorization as hypothesized. We will induct on m . If $m = 0$, then all the q_\bullet multiply out to 1 and hence by units, which makes no sense if $n > 0$. As such, the right-hand side must also be empty, meaning that $n = 0$, so there is nothing to prove. Note that a symmetric argument deduces that $n = 0$ implies $m = 0$, so we may assume that $m, n > 0$ in the argument which follows.

Now, for the induction, the hypothesis tells us that the irreducible p_m is prime and therefore must divide some factor q_\bullet on the right-hand side. Adjusting via a permutation, we may assume that $p_m \mid q_n$. Then we may write $q_n = u p_m$ for some $u \in A$, but in fact because p_m fails to be a unit, we see u is a unit because q_n

is irreducible. As such, adjusting by a unit, we may assume that $q_n = p_m$, whereupon dividing both of our factorizations out by this redundancy leaves us with

$$\prod_{i=1}^{m-1} p_i = \prod_{j=1}^{n-1} q_j,$$

and now we may induct downwards. ■

In fact, one has the following converse to Lemma A.14.

Lemma A.15. If A is a unique factorization domain, then an element $p \in A$ is prime if and only if it is irreducible.

Proof. The forward direction is covered by Lemma A.7. For the converse, suppose p is irreducible. Certainly p is not zero and not a unit. Now, suppose $p \mid ab$ for some $a, b \in A$. If $a = 0$, then $p \mid a$; similar holds if $b = 0$. Otherwise, we may give a and b factorizations into irreducibles by

$$a = \prod_{i=1}^m p_i \quad \text{and} \quad b = \prod_{j=1}^n q_j.$$

Because p divides the product ab , we see that $ab/p \in A$ will also have a factorization

$$\prod_{k=1}^s r_k = \frac{ab}{p},$$

so

$$p \prod_{k=1}^s r_k = ab = \prod_{i=1}^m p_i \cdot \prod_{j=1}^n q_j.$$

By the uniqueness of our factorizations, we see that the irreducible p (perhaps times a unit) must appear as one of the p_\bullet or q_\bullet , so $p \mid a$ or $p \mid b$. ■

We are now ready to prove Theorem A.11.

Proof of Theorem A.11. By Lemmas A.12 and A.14, it suffices to do the following two checks.

- Suppose we have an ascending chain of principal ideals

$$(a_0) \subseteq (a_1) \subseteq (a_2) \subseteq \cdots,$$

and we want to show that it stabilizes. Well, the union

$$I := \bigcup_{i=0}^{\infty} (a_i) = (a_0, a_1, a_2, \dots)$$

is an ideal of A . But all ideals are principal, so we may write $I = (a)$ for some $a \in A$. But then $a \in (a_N)$ for some N , so $I \subseteq (a_N)$, so for any $n \geq N$, we see that

$$(a_n) \subseteq I \subseteq (a_N) \subseteq (a_n),$$

establishing that our chain of ideals has stabilized.

- Note that an element of A is prime if and only if it is irreducible by combining Lemma A.7 and proposition A.10. ■

A.2 A Little Field Theory

The notes assume ring and group theory, but we will spend this appendix establishing the field theory that we will need.

A.2.1 Basic Notions

Here is our following definition.

Definition A.16 (field). A field K is a ring where each nonzero element has a multiplicative inverse.

We will be interested in how fields relate to each other.

Definition A.17 (field extension). A field extension L/K is when one field K is contained in another L . The degree $[L : K]$ of the extension is the dimension of L as a K -vector space. The field extension is said to be *finite* if and only if $[L : K] < \infty$.

Example A.18. Fix a field extension L/K . Given $\alpha \in L$, we can construct the field $K(\alpha)$ which is the quotient field of the integral domain generated by K and $\alpha \in L$. Formally, $K(\alpha)$ is the quotient field of the ring $K[\alpha]$ which is defined as the image of the ring homomorphism

$$\text{ev}_\alpha: K[x] \rightarrow L$$

defined by sending the polynomial $f(x) \in K[x]$ to $f(\alpha)$.

It might feel a little weird that we have jumped directly to one field being contained in another instead of a tamer notion of homomorphism. The following result explains why.

Lemma A.19. Let $\varphi: K \rightarrow L$ be a ring homomorphism of fields. Then φ is injective.

Proof. It suffices to check that $\ker \varphi$ is trivial. Well, $\ker \varphi$ is an ideal of K , but if nontrivial, then $\ker \varphi$ contains a unit and therefore must be all of K . However, $\varphi(1) = 1$, so $\ker \varphi \neq K$, so we see that we must have $\ker \varphi = 0$. ■

A basic fact about our extensions is how degrees behave in extensions.

Lemma A.20. Let M/L and L/K be extensions of fields. Then

$$[M : L][L : K] = [M : K].$$

Proof. If M/L or L/K are infinite, then the K -vector space M will have an infinitely linearly independent set, meaning $[M : K]$ is also infinite. Otherwise, we take $[M : L] = [L : K]$ to be finite. Let $\{m_1, \dots, m_r\}$ and $\{\ell_1, \dots, \ell_s\}$ be bases for M/L and L/K , respectively. We claim that

$$\{m_i \ell_j\}_{1 \leq i \leq r, 1 \leq j \leq s}$$

is a basis for M/K .

- We show that the $m_i \ell_j$ span: any $m \in M$ can be expressed as

$$m = \sum_{i=1}^r a_i m_i$$

where $a_k \in L$, but then each $a_k \in L$ can be expressed as

$$m = \sum_{i=1}^r \sum_{j=1}^s b_{ij} m_i \ell_j,$$

where $b_{ij} \in K$. This is what we wanted.

- We show that the $m_i \ell_j$ are linearly independent: suppose

$$\sum_{i=1}^r \sum_{j=1}^s b_{ij} \ell_j m_i = \sum_{i=1}^r \sum_{j=1}^s b'_{ij} \ell_j m_i.$$

Then

$$\sum_{i=1}^r \sum_{j=1}^s (b_{ij} - b'_{ij}) \ell_j m_i = 0,$$

so because $\{m_1, \dots, m_r\}$ is a basis of M/L , we see

$$\sum_{j=1}^s (b_{ij} - b'_{ij}) \ell_j = 0$$

for each i , but because $\{\ell_1, \dots, \ell_s\}$ is a basis of L/K , we see $b_{ij} = b'_{ij}$ for each i and j . ■

A.2.2 Polynomial Rings

In this subsection, we show that $K[x]$ is a unique factorization domain for any field K . We will not bother to show the usual facts about degree over integral domains, such as $\deg fg = \deg f + \deg g$ for nonzero $f, g \in K[x]$. However, we will show the following result, whose proof is more technical than one would like. Thankfully, we will not have to use this result for a while.

Lemma A.21. Fix an irreducible polynomial $f \in \mathbb{C}[x]$. Then f has no repeated roots.

Proof. Note $\deg f \geq 1$ because $\deg f = 1$ implies that f is a unit and thus not irreducible. Because \mathbb{C} is algebraically closed, we note that $f(x)$ factors as

$$f(x) = c \prod_{i=1}^n (x - \alpha_i)$$

for some complex numbers $c, \alpha_1, \dots, \alpha_n \in \mathbb{C}$. Now, if $\alpha_i = \alpha_j$ for $i \neq j$, then $f(x)$ has a double root at α_i , so a direct computation shows that $f'(\alpha_i) = 0$, so $f(x)$ and $f'(x)$ have a root in common, so $\gcd(f(x), f'(x))$ is a non-constant polynomial with degree strictly less than $f(x)$ but dividing $f(x)$, which contradicts $f(x)$ being irreducible. ■

Anyway, the main point to showing that $K[x]$ is a unique factorization domain is the following result.

Proposition A.22 (division). Fix a field K , and let $a, b \in K[x]$ be polynomials with $b \neq 0$. Then there exist polynomials $q, r \in K[x]$ such that

$$a = bq + r$$

where $r = 0$ or $0 \leq \deg r < \deg b$.

Proof. We induct on $\deg a$. If $a = 0$ or $\deg a < \deg b$, we set $q = 0$ and $r = a$. Otherwise, $\deg a \geq \deg b$. Let the leading coefficient of b be $b_n x^n$, and let the leading coefficient of a be $a_m x^m$. Then

$$c(x) := a(x) - \frac{a_m}{b_n} x^{m-n} \cdot b$$

cancels out the leading coefficient of a , so $c = 0$ or $\deg c < \deg a$. So we can apply the result to c by the induction, writing

$$c = bq_c + r_c,$$

so

$$a(x) = b(x) \left(q_c(x) + \frac{a_m}{b_n} x^{m-n} \right) + r(x),$$

finishing. ■

Theorem A.23. Fix a field K . Then the field $K[x]$ is a principal ideal domain and hence a unique factorization domain.

Proof. If we show that $K[x]$ is a principal ideal domain, we finish immediately by Theorem A.11. So we want to show that $K[x]$ is a principal ideal domain.

Well, let $I \subseteq K[x]$ be a principal ideal domain. If $I = \{0\}$, then $I = (0)$, so there is nothing to say. Otherwise, I has a nonzero element, so we let $f \in I$ denote any element of least degree. We claim $I = (f)$. Certainly $(f) \subseteq I$, so we want to show that $a \in I$ lives in (f) . Well, by Proposition A.22, we may write

$$a = fq + r$$

where $r = 0$ or $0 \leq \deg r < \deg f$. However, $r = a - fq \in I$ has $\deg r < \deg f$, so minimality of $\deg f$ requires $r = 0$, meaning $a = fq$, so $a \in (f)$. ■

A.2.3 Algebraic Elements

In this subsection, we show some basic properties of algebraic elements. Here is our definition.

Definition A.24 (algebraic). Fix a field extension L/K . An element $\alpha \in L$ is *algebraic over K* if and only if α is the root of some nonzero polynomial in $K[x]$.

Example A.25. Any element of K is algebraic over K because $\alpha \in K$ is the root of the polynomial $x - \alpha \in K[x]$.

Finite extensions provide a wealth of algebraic elements.

Lemma A.26. Let L/K be a finite extension of fields. Then each $\alpha \in L$ is algebraic over K .

Proof. The elements $1, \alpha, \alpha^2, \dots$ form an infinite set in L , so they cannot be K -linearly independent because $\dim_K L < \infty$. Thus, there is a relation of the form

$$\sum_{k=0}^n a_k \alpha^k = 0$$

where $a_k \in K$ are not all zero. As such, the polynomial $f(x) := \sum_{k=0}^n a_k x^k$ will do. ■

It will shortly be helpful to limit the polynomial attached to α somewhat.

Lemma A.27. Fix a field extension L/K , and let $\alpha \in L$ be algebraic over K . Then α is the root of a unique monic irreducible polynomial $f(x) \in K[x]$. In fact, for any polynomial $g \in K[x]$ with $g(\alpha) = 0$, we have $f \mid g$.

Proof. We begin by showing existence. We know that α is the root of some nonzero polynomial $f(x) \in K[x]$, so we choose $f(x)$ to have the smallest degree possible. By dividing out the leading coefficient (which is nonzero because f is nonzero), we may assume that f is monic. It remains to show that f is irreducible. Well, suppose that

$$f = ab$$

for $a, b \in K[x]$. Note neither a nor b is zero because this would imply $f = 0$; additionally, if both are units, then f is a unit and hence a constant polynomial, which also makes no sense. Now, evaluating at α , we see that $f(\alpha) = 0$ requires $a(\alpha) = 0$ or $b(\alpha) = 0$, so by minimality of f , we must have $\deg a \geq \deg f$ or $\deg b \geq \deg f$. Without loss of generality take $\deg a \geq \deg f$, but $f = ab$ then forces $\deg a = \deg f$, and b is a constant polynomial.

We now show that $g(\alpha) = 0$ implies $f \mid g$ for any $g \in K[x]$. By Proposition A.22, we may write

$$g = fq + r$$

where $r = 0$ or $0 \leq \deg r < \deg f$. By plugging in α , we see that $f(\alpha) = g(\alpha) = 0$ implies $r(\alpha) = 0$. But if nonzero $\deg r < \deg f$, violating minimality of f , so we instead have $r = 0$, implying $g = fq$ and so $f \mid g$.

To finish up, we show that f is unique. Well, if g is another monic irreducible polynomial with $g(\alpha) = 0$, then $f \mid g$ by the above argument. But f is nonzero, so $f \mid g$ requires $g = fu$ for a unit u . Being a unit means that u is a constant polynomial in K , so for example $\deg f = \deg g$, and because f and g have the same leading coefficient, we must have $u = 1$. Thus, $f = g$, as needed. ■

Lemma A.28. Fix a field extension L/K , and let $\alpha \in L$ be algebraic over K and in particular the root of a monic irreducible polynomial $f(x) \in K[x]$. Then

$$\frac{K[x]}{(f(x))} \cong K[\alpha].$$

In particular, $K[\alpha]$ is a field of degree $\deg f$ over K .

Proof. Quickly, note that the last sentence follows from the isomorphism because

For the first claim, note that there is a surjective ring homomorphism $\text{ev}_\alpha: K[x] \rightarrow K[\alpha]$ by sending $g(x) \mapsto g(\alpha)$ for any $g(x) \in K[x]$. We want to show that $\ker \text{ev}_\alpha = (f)$. Certainly $f \in \ker \text{ev}_\alpha$. For the other inclusion, we note that any $g \in \ker \text{ev}_\alpha$ has $g(\alpha) = 0$ and hence $f \mid g$ by Lemma A.27.

For the second claim, note that $K[\alpha]$ is not a field because $K[x]/(f(x))$ is a field by Remark A.9. As for the degree computation, write

$$f(x) = x^d + \sum_{k=0}^{d-1} a_k x^k.$$

Then for each $n \geq d$, we can express α^n in terms of α^k with $k < n$: indeed, $f(\alpha) = 0$ implies

$$\alpha^n = - \sum_{k=0}^{d-1} a_k \alpha^{k+n-d}.$$

Thus, $1, \alpha, \alpha^2, \dots, \alpha^{d-1}$ spans $K[\alpha]$, so $\dim_K K[\alpha] \leq d$. In fact, these α^k are linearly independent because any nontrivial relation involving them becomes a polynomial $g(x)$ with α a root which is either zero or has degree less than $\deg f$, but Lemma A.27 enforces $g = 0$. ■

As a last aside, we note that the sum and product of algebraic elements remains algebraic. This requires a trick known as the “determinant trick.”

Proposition A.29. Fix a field extension L/K and some $\alpha \in L$. Then the following are equivalent.

- (a) α is algebraic over K .
- (b) The field $K[\alpha]$ is a finite extension of K .
- (c) There is a subfield $K' \subseteq L$ finite over K which contains α .

Proof. We show the implications separately.

- Here, (a) implies (b) is proven in Lemma A.28.
- Note (b) implies (c) by setting $K' := K[\alpha]$.
- Checking that (c) implies (a) is harder. Suppose K' is generated by the elements $\alpha'_1, \alpha'_2, \dots, \alpha'_n$. Note that $\alpha\alpha'_i \in A'$ for each α'_i , so we may write

$$\alpha\alpha'_i = \sum_{j=1}^n a_{ij}\alpha'_j$$

for some elements $a_{ij} \in K'$. In other words, the matrix $T := (a_{ij})_{i,j=1}^n$ has

$$\alpha \begin{bmatrix} \alpha'_1 \\ \vdots \\ \alpha'_n \end{bmatrix} = T \begin{bmatrix} \alpha'_1 \\ \vdots \\ \alpha'_n \end{bmatrix}.$$

Thus, $T - \alpha I_n$ is an $n \times n$ matrix with entries in K' , and it has the nonzero vector $(\alpha'_1, \dots, \alpha'_n)$ in its kernel, so $\det(T - \alpha I_n) = 0$. Expanding out the polynomial $\det(\alpha I_n - T) = 0$ makes α the root of a monic polynomial (of degree n) with coefficients in A , so α is indeed integral over A . ■

Corollary A.30. Fix a field extension L/K , and let K' denote the set elements of L algebraic over K . Then K' is a subfield of L . In fact, for any $\alpha \in L$ algebraic over K' , we have $\alpha \in K'$.

Proof. We run our checks separately.

- We check that K' is a field. Note $0, 1 \in K'$ because these elements are the roots of the polynomials x and $x-1$, respectively. It remains to show that, for any $\alpha, \beta \in K'$, we have $\alpha+\beta, \alpha\beta \in K'$ and $\alpha/\beta \in K'$ if $\beta \neq 0$. The main point is to show that $K[\alpha, \beta]$ is a finite extension of K , which will complete the proof by Proposition A.29.

Well, let α and β be the roots of the monic polynomials $f(x), g(x) \in K[x]$ respectively. Then by Proposition A.29 shows that $K[\beta]$ is a finite extension of K , and $f(\alpha) = 0$ shows that α is integral over $K[\beta]$, so $K[\alpha, \beta]$ is finite field extension of $K[\beta]$. We conclude $K[\alpha, \beta]$ is a finite field extension of K by Lemma A.20.

- Suppose that $\alpha \in L$ is the root of the monic polynomial $f(x) \in K'[x]$ (monic by Lemma A.27); we show that $\alpha \in K'$. Well, expand $f(x)$ as

$$f(x) = x^d + \sum_{k=0}^{d-1} a_k x^k$$

for some $a_0, \dots, a_{d-1} \in K'$. Each a_\bullet is algebraic over K , so Proposition A.29 tells us that $K[a_\bullet]$ for each a_\bullet . As such, as in the previous check, we may build the tower

$$K \subseteq K[a_0] \subseteq K[a_0, a_1] \subseteq \dots \subseteq K[a_0, \dots, a_{d-1}],$$

where each field is finite over the previous one by Proposition A.29. Then Lemma A.20 tells us that $K[a_0, \dots, a_{d-1}]$ is finite over K . Lastly, $f(\alpha) = 0$ tells us that α is algebraic over $K[a_0, \dots, a_{d-1}]$, so $\mathbb{Z}[a_0, \dots, a_{d-1}, \alpha]$ is finite over $K[a_0, \dots, a_{d-1}]$ —and hence finite over K by Lemma A.20, meaning that α is algebraic over K by Proposition A.29. ■

A.2.4 Enough Galois Theory to be Dangerous

We are going to derive a lot of mileage from the following result in field theory. It leads towards Galois theory; even though Galois theory is a beautiful subject, it is one that we can avoid somewhat.

Proposition A.31. Let L/K be a finite field extension, where L is a subfield of \mathbb{C} . Then each embedding $\sigma: K \rightarrow \mathbb{C}$ extends to exactly $[L : K]$ embeddings $\tilde{\sigma}: L \hookrightarrow \mathbb{C}$.

Here, we are using the term “embedding” to refer to an (injective) ring homomorphism.

Proof. We induct on $[L : K]$, which is legal because $[L : K] < \infty$. If $[L : K] = 1$, then $L = K$, and there is nothing to say because we must have $\tilde{\sigma} = \sigma$.

Otherwise, suppose $[L : K] > 1$. Then fix $\alpha \in L \setminus K$. By Lemma A.26, α is algebraic over K , so by Lemma A.27, α is the root of some monic irreducible polynomial $f(x)$. Now, \mathbb{C} is algebraically closed, so we note that $f(x)$ factors as

$$f(x) = \prod_{i=1}^n (x - \alpha_i)$$

for some complex numbers $\alpha_1, \dots, \alpha_n \in \mathbb{C}$; note that the α_i are distinct by Lemma A.21. Thus, given an embedding $\sigma: K \hookrightarrow \mathbb{C}$, there are n exactly extensions to $\sigma_i: K[\alpha] \hookrightarrow \mathbb{C}$ by sending $\sigma_i(\alpha) := \alpha_i$. We have a number of checks to make this sentence make sense.

- Setting $\sigma_i(\alpha) = \alpha_i$ defines a unique embedding $K[\alpha] \hookrightarrow \mathbb{C}$. The embedding here is uniquely defined because we need to have $\sigma_i|_K = \sigma$, and then any polynomial in $K[\alpha]$ will have its output determined by where α goes. To show that σ_i is well-defined, we note that it is simply the composite

$$K[\alpha] \cong \frac{K[x]}{(f(x))} \cong K[\alpha_i] \subseteq \mathbb{C},$$

where the left isomorphism is by Lemma A.28.

- We have in fact defined n embeddings because the roots α_i are distinct.
- Each extension $\tilde{\sigma}: K[\alpha] \rightarrow \mathbb{C}$ of σ must take this form. It suffices by our first point to check that $\tilde{\sigma}(\alpha) = \alpha_i$ for some α_i . Well, note that

$$f(\tilde{\sigma}(\alpha)) = \tilde{\sigma}(f(\alpha)) = 0$$

because $\tilde{\sigma}$ is a ring homomorphism. The result follows.

Now, by induction each of the σ_i extend to exactly $[L : K[\alpha]] < [L : K]$ distinct embeddings $L \hookrightarrow \mathbb{C}$, totaling to

$$[L : K[\alpha]] \cdot [K[\alpha] : K] = [L : K]$$

embeddings $L \hookrightarrow \mathbb{C}$, where we have used Lemma A.20. Let σ_i extend to $\sigma_{i1}, \dots, \sigma_{im}$ where $m = [L : K[\alpha]]$. We have the following checks on the σ_{ij} .

- Note that σ_{ij} must be distinct: if $\sigma_{ij} = \sigma_{i'j'}$, then restricting to $K[\alpha]$ reveals that $\sigma_{ij}|_{K[\alpha]} = \sigma_{i'j'}|_{K[\alpha]} = \sigma_i$, so $\sigma_i = \sigma_{i'}$, so $i = i'$. But then the uniqueness of extending from $K[\alpha]$ to L means that $\sigma_{ij} = \sigma_{i'j'}$ implies $j = j'$.
- We show that all extensions $\tilde{\sigma}: L \hookrightarrow \mathbb{C}$ of σ take the form σ_{ij} . Well, restricting $\tilde{\sigma}$ to $K[\alpha]$ shows that $\tilde{\sigma}|_{K[\alpha]} = \sigma_i$ for some i . Then $\tilde{\sigma} = \sigma_{ij}$ for some j by construction of the σ_{ij} .

The above checks show that the σ_{ij} provide all extensions of $\sigma: K \hookrightarrow \mathbb{C}$, counted uniquely, finishing. ■

A.2.5 Norm and Trace

Proposition A.31 allows us to make sense of the norm and trace of an algebraic element, which we now define.

Definition A.32. Let L/K be a finite extension of fields. Then for $\alpha \in L$, let $\mu_\alpha: L \rightarrow L$ denote the multiplication-by- α map, which is K -linear by the distributive law. Then we define the *trace* of α as $T_{L/K}(\alpha) := \text{tr } \mu_\alpha$ and the *norm* of α as $N_{L/K}(\alpha) := \det \mu_\alpha$.

Example A.33. Let L/K be a finite extension of fields. Then any $\alpha \in K$ has μ_α given by the matrix $\alpha I_{[L:K]}$, so $T_{L/K}(\alpha) = [L:K]\alpha$ and $N_{L/K}(\alpha) = \alpha^{[L:K]}$.

Here is our key example of the norm and trace.

Proposition A.34. Let L/K be an extension of fields. Then let $\alpha \in K$ be algebraic over K which is the root of the monic irreducible polynomial $f(x) \in K[x]$. Writing $f(x) = x^d + \sum_{k=0}^{d-1} a_k x^k$, we have

$$T_{K[\alpha]/K}(\alpha) = -a_{d-1} \quad \text{and} \quad N_{K[\alpha]/K}(\alpha) = (-1)^d a_0.$$

Proof. Note $K[\alpha]$ is finite over K by Lemma A.28, where we actually showed that $1, \alpha, \alpha^2, \dots, \alpha^{d-1}$ is a basis of L as a K -vector space. Then μ_α with respect to this (ordered) basis looks like the $d \times d$ matrix

$$\begin{bmatrix} 0 & 0 & \cdots & 0 & -a_0 \\ 1 & 0 & \cdots & 0 & -a_1 \\ 0 & 1 & \cdots & 0 & -a_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & -a_{d-1} \end{bmatrix}. \quad (\text{A.1})$$

The trace of this matrix is $-a_{d-1}$, and its determinant is $(-1)^d a_0$ by expansion by minors. ■

Corollary A.35. Let K/\mathbb{Q} be a finite extension of fields of degree n . Fix $\alpha \in K$, and let $\sigma_1, \dots, \sigma_n$ denote the embeddings $K \hookrightarrow \mathbb{C}$. Then

$$T_{K/\mathbb{Q}}(\alpha) = \sum_{i=1}^n \sigma_i(\alpha) \quad \text{and} \quad N_{K/\mathbb{Q}}(\alpha) = \prod_{i=1}^n \sigma_i(\alpha).$$

Proof. We use Proposition A.34. Note $\alpha \in K$ is algebraic over \mathbb{Q} by Lemma A.26, so let τ_1, \dots, τ_d denote the embeddings $K[\alpha] \hookrightarrow \mathbb{C}$. By the proof of Proposition A.31, we see that

$$f(x) = \prod_{i=1}^d (x - \tau_i(\alpha)),$$

where $f(x) \in \mathbb{Q}[x]$ is the unique monic irreducible polynomial with $f(\alpha) = 0$ provided by Lemma A.27. Thus, Proposition A.34 tells us that

$$T_{\mathbb{Q}[\alpha]/\mathbb{Q}}(\alpha) = \sum_{i=1}^d \tau_i(\alpha) \quad \text{and} \quad N_{\mathbb{Q}[\alpha]/\mathbb{Q}}(\alpha) = \prod_{i=1}^d \tau_i(\alpha).$$

To complete the proof, we must extend up from $K[\alpha]/K$ to L/K . Well, let $\ell_1, \dots, \ell_{n/d}$ denote a basis for L as a $K[\alpha]$ -vector space, where we are implicitly using Lemma A.20. Then the proof of Lemma A.20 shows us that $\ell_i \alpha^j$ provides a basis for L/\mathbb{Q} , so writing out μ_α according to this basis looks like n/d blocks of (A.1). Because each τ_i extends to exactly n/d embeddings $K \hookrightarrow \mathbb{C}$ by Proposition A.31, the result follows. ■

BIBLIOGRAPHY

- [Old70] C. D. Olds. “The Simple Continued Fraction Expansion of e ”. In: *The American Mathematical Monthly* 77.9 (1970), pp. 968–974. ISSN: 00029890, 19300972. URL: <http://www.jstor.org/stable/2318113> (visited on 08/26/2023).
- [HW75] G. H. Hardy and E. M. Wright. *An Introduction to the Theory of Numbers*. Oxford, 1975.
- [Kle16] Felix Klein. *Elementary Mathematics from a Higher Standpoint*. Trans. by Gert Schubring. Vol. II. Springer Berlin, Heidelberg, 2016.
- [Shu16] Neal Shusterman. *Scythe*. Arc of a Scythe. Simon & Schuster, 2016.
- [Pro22] Ross Mathematics Program. *Students*. 2022. URL: <https://rossprogram.org/students/>.
- [Con] Keith Conrad. *Transcendence of e* . URL: <https://kconrad.math.uconn.edu/blurbs/analysis/transcendence-e.pdf>.

LIST OF DEFINITIONS

algebraic, [37](#), [78](#)
algebraic integer, [64](#)

continued fraction, [11](#)
convergent, [14](#), [24](#)

field, [76](#)
field extension, [76](#)

infinite continued fraction, [21](#)
integral, [65](#)
integral domain, [72](#)

irrationality measure, [32](#)
irreducible, [72](#)

normal, [64](#)
number ring, [65](#)

prime, [72](#)

transcendental, [37](#)

unique factorization domain, [72](#)