

DIOPHANTINE APPROXIMATION AND DYNAMICAL SYSTEMS

ANTONY LEE

An exploration of rotation sequences and the farey tessellation

Lunds Tekniska Högskola

February 2020 – version 1.2

FOREWORD

This project was written at Lunds Tekniska Högskola as part of a degree project for the technical mathematics program. The work took place throughout most of 2019, and it was finalised early in 2020.

The goal was to write an educational text describing the connection between three seemingly disparate areas of mathematics. As such the work has been a balance between picking out the relevant parts from each subject, understanding each of them to a high enough degree as to be able to explain them, and judging which parts might need more exposition, and what I can assume the reader will already know. I hope that these relatively diverse, yet connected fragments of mathematics will convey some of the moments of wonder I had while learning about them.

Due to the width of the knowledge required this work would most likely have been infeasible without the help of my supervisor, Jörg Schmeling, whom I would like to thank for his great amount of help and patience in advising me, as well as for suggesting the subject to me in the first place.

I would also like to thank my family for their love and support. I do not know what I would have done without them.

ABSTRACT

The subject of diophantine approximation is a classical mathematic problem, as old as it is well studied. There are many different texts describing its connection to more modern areas of study, but few which do so with the aim of exploring the connections themselves. This paper aims to serve as an introduction to diophantine approximation, and to expose some properties common between two dynamical systems where it occurs. This is done in the style of a booklet, starting from the basics in each of the areas of diophantine approximation, continued fractions, symbolic sequences, and hyperbolic geometry. Focus on each of the chapters following the first is on how to they connect back to diophantine equation. The chapters are then capped off with additional notes which explore things related to their respective subjects, for example the modern advancements made in the subject, or other interesting trivia for the interested reader.

For complete comprehension of the text, the reader is assumed to have basic knowledge of the relation between rational and real numbers, analysis, matrices, number theory, and function theory.

The text largely succeeds in its goals as an educatory text and is thought to be a somewhat novel contribution to the body of literature on the subject. Further work could expand on this by incorporating further areas in mathematics where diophantine approximation appears. Another avenue of exploration is to explore the underlying reasons for the similarities exposed here.

ABSTRAKT (SVENSKA)

Diofantisk approximering är ett klassiskt problem inom matematiken som är lika gammalt som det är studerat. Det finns många olika texter som kopplar samman området till mer moderna idéer, men få som undersöker samband i hur själva kopplingen förkroppsligas. Denna artikeln har som syfte att leda läsaren genom kunskapen som behövs för att förstå särskilda egenskaper gemensamma till två olika dynamiska system, för att sedan visa upp dessa egenskaper. Detta görs i form av en handbok som i tur och ordning går igenom grunderna i diofantisk approximering, kedjebråk, symboliska följder, samt hyperbolisk geometri, där det för var och en av de tre sista förtydligas hur de hänger ihop med det första. Varje avsnitt avslutas därefter med mindre utflykter till ting relaterat till området i fråga, exempelvis moderna utvecklingar inom området eller annan information som kan vara till intresse för läsaren.

För att fullt ut kunna tillgodogöra sig texten antas läsaren ha åtminstone grundläggande kunskap inom ett antal områden, bland annat förhållandet mellan rationella och rella tal, endimensionell analys, matriser, talteori, samt även funktionsteori.

Texten fullföljer i stort sett sitt mål att vara en resurs för intresserade läsare att lära sig om området, och är till författarens vetskap även innehållsmässigt ett nytt bidrag till mängden litteratur om ämnet. Ytterligare arbete kan göras genom att utforska andra områden där diofantisk approximering visar sig. Alternativt skulle orsaken till sambanden som uppvisas i detta arbetet kunna utforskas på en djupare nivå.

CONTENTS

1	FOREWORD	iii
2	DIOPHANTINE APPROXIMATIONS	1
3	CONTINUED FRACTIONS	7
4	SYMBOLIC SEQUENCES	23
5	HYPERBOLIC GEOMETRY	31
	BIBLIOGRAPHY	45

DIOPHANTINE APPROXIMATIONS

With historical roots stretching back to at least the third century, Diophantine approximation is the approximation of the reals using rational numbers. It is named after Diophantus of Alexandria, and is closely related to Diophantine equations, where one tries to find all natural number solutions to some given equation.

Some numbers are easier to approximate by rationals than others. Trivially, any rational number can be perfectly approximated (by itself), while irrational numbers cannot. There is however some variation in approximability among the irrational numbers. Take for example Dirichlet's approximation theorem.

Theorem 2.0.1 (Dirichlet). *For any $\alpha \in \mathbb{R}$, and $N \in \mathbb{N}$ there exist integers p and q such that*

$$\left| \alpha - \frac{p}{q} \right| < \frac{1}{qN} \leq \frac{1}{q^2}, 1 \leq q \leq N.$$

Proof. This can be proven by the pigeonhole principle. Denote the integer part of a real number x by $[x]$, and the fractional part by $\{x\}$. Now consider $\{n\alpha\}$ for $n = 0, 1, 2, \dots, N$, the $N + 1$ pigeons. The holes are the N non-overlapping intervals of length $\frac{1}{N}$, $[\frac{m}{N}, \frac{m+1}{N})$, where $m = 0, 1, 2, \dots, N - 1$. Since the intervals completely cover $[0, 1)$, there must be at least one interval where two different $\{n\alpha\}$ lie. Call them $\{n_1\alpha\}$ and $\{n_2\alpha\}$, such that $n_1 < n_2$. This means

$$|\{n_2\alpha\} - \{n_1\alpha\}| < \frac{1}{N}.$$

Recall $x = [x] + \{x\}$, so $\{n\alpha\} = n\alpha - [n\alpha]$, and

$$|(n_2 - n_1)\alpha - ([n_2\alpha] - [n_1\alpha])| < \frac{1}{N}.$$

Identify $p = [n_2\alpha] - [n_1\alpha]$, and $q = n_2 - n_1$, which finally gives

$$\left| \alpha - \frac{p}{q} \right| < \frac{1}{qN} \leq \frac{1}{q^2}.$$

□

Since $|\alpha - \frac{p}{q}|$ is to be minimised, we can also rewrite it into another almost equivalent form, the exact difference being something we will come back to later. Begin by multiplying by q , thus eliminating the fraction, giving $|q\alpha - p|$. To minimise the expression we always choose p to be the integer nearest to $q\alpha$. It then makes sense to

introduce the notation $\|x\|$, the distance from x to the nearest integer, and we end up with $\|q\alpha\|$. This form gives some additional insight, as apparently the numerator p could be eliminated. In this sense, only the choice of the denominator q is important when making approximations by rationals. Dirichlet's theorem can be restated with this form as

Theorem 2.0.2 (Dirichlet). *For any $\alpha \in \mathbb{R}$, and $N \in \mathbb{N}$, there exists a natural number $1 \leq q \leq N$ such that*

$$\|q\alpha\| < \frac{1}{N} \leq \frac{1}{q}.$$

The same proof as before still applies. Dirichlet's theorem is essentially about the approximability of irrational numbers, here named α , and gives us an upper bound for the worst possible approximation of any given q . Although the bound does improve as q grows, an increase in q does not necessarily yield a better approximation.

Definition 2.0.3. A *best* rational approximation of an irrational number, α , is a fraction, $\frac{p^*}{q^*}$, such that for all other fractions $\frac{p}{q}$, where $0 < q \leq q^*$,

$$\left| \alpha - \frac{p^*}{q^*} \right| < \left| \alpha - \frac{p}{q} \right|.$$

Since it is possible to get arbitrarily close to α by increasing q , there are infinitely many best rational approximations for any irrational α . Take for example when $\alpha = \pi$, for which $\frac{22}{7} \approx 3.142857$ is often given as a best approximation. By instead looking at $q = 8$, we find the two closest approximations to be $\frac{25}{8} = 3.125$ and $\frac{26}{8} = 3.25$, both of which are further away from π than $\frac{22}{7}$. In other words, $q = 8$ does not give an improved approximation over $q = 7$ for any choice of p .

Let us continue by defining the most difficult numbers to approximate.

Definition 2.0.4. A number, α , is called *badly approximable* if

$$\left| \alpha - \frac{p}{q} \right| > \frac{c}{q^2},$$

or alternatively,

$$\|q\alpha\| > \frac{c}{q},$$

for some constant $c > 0$ and all rational numbers $\frac{p}{q}$.

This will be expanded slightly upon later, but *badly approximable* can be thought of as *slowly converging* in the sense that large denominators are needed to approximate the number accurately.

Similarly, there is a measure for the irrationality of a number.

Definition 2.0.5. Let R be the set of positive reals μ for which $\|q\alpha\| < \frac{1}{q^\mu}$ has finitely many solutions. The *irrationality measure* of an irrational real number α , denoted $\mu(\alpha)$, is defined by

$$\mu(\alpha) := \inf_{\mu \in R} \mu.$$

It can be shown, for example, that rational numbers have $\mu = 1$, and transcendental numbers have $\mu \geq 2$.

There are some numbers for which R is empty, and so for the sake of consistency, let us define $\mu(\alpha) = \infty$ for such numbers, which leads us to the following definition.

Definition 2.0.6. A *Liouville number* is a number α for which $0 < \|q\alpha\| < \frac{1}{q^\mu}$ has infinitely many integer solutions q for all μ , or equivalently, a number for which $\mu(\alpha) = \infty$.

Liouville numbers are irrational numbers that can be approximated very closely by rational numbers. As a corollary from the below theorem, all irrational Liouville numbers are transcendental.

Theorem 2.0.7 (Liouville). *If α is an irrational algebraic number of degree $n \geq 2$ (i.e. irrational roots of polynomials of order 2 or higher), then there exists a real constant $c > 0$ such that*

$$\left| \alpha - \frac{p}{q} \right| > \frac{c}{q^n},$$

for all $\frac{p}{q}$.

Proof. Let $f \in \mathbb{Z}[x]$ be a non-constant irreducible polynomial of degree n , such that $f(\alpha) = 0$. We may assume integer coefficients, since any fractions can be eliminated without changing the roots, by multiplying the polynomial by an appropriate number.

According to the mean-value theorem, for a given $\frac{p}{q}$, there exists a ξ such that

$$f'(\xi) = \frac{f(\alpha) - f\left(\frac{p}{q}\right)}{\alpha - \frac{p}{q}} = -\frac{f\left(\frac{p}{q}\right)}{\alpha - \frac{p}{q}}.$$

Taking absolute values, we find that

$$|f'(\xi)| = \frac{|f\left(\frac{p}{q}\right)|}{\left|\alpha - \frac{p}{q}\right|}.$$

After some rearranging, we get

$$\left| \alpha - \frac{p}{q} \right| = \left| f\left(\frac{p}{q}\right) \right| \cdot |1/f'(\xi)| \geq \frac{|1/f'(\xi)|}{q^n}.$$

The last inequality follows if you note that f only has integer coefficients and is of degree n , so inserting the fraction $\frac{p}{q}$ must give a

fraction with a denominator that is at most q^n . Since f is irreducible, there must be no rational (integral) roots, so the numerator cannot be zero.

As ξ was chosen from the mean value theorem, it lies inbetween $\frac{p}{q}$ and α . Then for sufficiently good approximations $\alpha \approx \frac{p}{q}$, we have $f'(\xi) \approx f'(\alpha)$, such that for example $|f'(\xi)| \geq \frac{1}{2} \cdot |f'(\alpha)|$. We find

$$\left| \alpha - \frac{p}{q} \right| \geq \frac{2/|f'(\alpha)|}{q^n}.$$

Due to irreducibility is that $f'(\alpha) \neq 0$, as otherwise α would be a double root, so $2/|f'(\alpha)|$ is well defined.

Then the result follows by letting $c = 2/|f'(\alpha)|$. Since the inequality $|f'(\xi)| \geq \frac{1}{2} \cdot |f'(\alpha)|$ was chosen rather arbitrarily, c can be lowered to allow for smaller q if necessary. \square

In actuality, the above theorem is not as strict as possible. A version that is optimal with respect to the right hand side was developed through a series of successive improvements, the final being due to K. Roth for which he was awarded a Fields medal.

Theorem 2.0.8 (Thue-Siegel-Roth). *If α is an algebraic number of degree $n \geq 2$, then for all $\epsilon > 0$,*

$$\left| \alpha - \frac{p}{q} \right| < \frac{1}{q^{2+\epsilon}}$$

has infinitely many integer solutions p and q .

Another way of stating the same thing is that all irrational algebraic numbers of degree 2 or higher has the irrationality measure $\mu = 2$. This means that irrational roots of polynomials are among the most difficult to approximate by rational numbers, which might be somewhat unexpected since they in other contexts might be considered to be "pleasant" numbers. Among them is the worst approximable number, the golden ratio, $\phi = \frac{1+\sqrt{5}}{2}$. Using ϕ as a worst-case scenario in Dirichlet's theorem, it is in fact possible to achieve a sharpening.

Theorem 2.0.9. $|\alpha - \frac{p}{q}| \leq \frac{1}{\sqrt{5}q^2}$

We shall leave it to the interested reader to find a full proof in [10], but a clue on the presence of ϕ can be seen in the $\sqrt{5}$ in the denominator.

2.0.0.1 Additional notes

V. Jarnik showed that the set of all badly approximable numbers (BANs) have Hausdorff dimension 1. The Hausdorff dimension can be seen as one way to measure the dimension of fractals embedded in metric spaces. The real line is not fractal in its nature, and so it retains

its Hausdorff dimension as $\mathbb{1}$. Since the BANs have the same Hausdorff dimension as the real numbers, one might conclude that nearly all numbers are badly approximable. An interesting counterpoint is that the BANs also happen to have a Lebesgue measure (typically used for integration) of 0 , suggesting that nearly no number is badly approximable. As such, the commonness of BANs depends on ones perspective.

A famous, and as of the time of writing, unsolved conjecture in the area of diophantine approximation is the Littlewood conjecture.

Conjecture 2.0.10 (Littlewood). *For any $\alpha, \beta \in \mathbb{R}$,*

$$\liminf_{n \rightarrow \infty} n \|n\alpha\| \|n\beta\| = 0.$$

Some of the meaning behind the conjecture can be gleaned by first thinking about the approximability of α and β . Most notably the conjecture is trivial if either is not a badly approximable number, as the limit then goes to zero. It is also clear that while we are only interested in the case when α and β are badly approximable, the statement could not hold if neither of the two respective norms goes to zero, due to the presence of the factor n . This means that we need good (enough) rational approximations of α and β . For example for some $S \subset \mathbb{N}$, we could assume that

$$\liminf_{m \rightarrow \infty} m \|m\alpha\| = C,$$

where $C > 0$ is a constant and $m \in S$. Notably, m does *not* need to be *best* approximations of α . From this we can see the things in terms of diophantine approximation, the new interpretation of the conjecture being that we can always find a set that approximates both α and β well (enough) simultaneously. In other words, that there always is a set

$$S' \in \left\{ S \mid \liminf_{x \rightarrow \infty} x \|x\alpha\| < \infty, x \in S \right\}$$

such that

$$\liminf_{m \rightarrow \infty} m \|m\beta\| = 0, m \in S'.$$

CONTINUED FRACTIONS

Deeply connected to diophantine approximation are the continued fractions. A *continued fraction* is a fraction-type expression of a certain form,

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

where $a_0 \in \mathbb{Z}$ and $a_k \in \mathbb{N}, k > 0$. Since (typically) only the a_k change from one continued fraction to another, they can instead be represented in a more compact form,

$$\alpha = [a_0; a_1, a_2, \dots, a_n].$$

One may wonder why the coefficients, a_k , are restricted to the naturals. There are some cases where it is appropriate to allow for real coefficients, so the term *regular* (alternatively *simple*) may then be used to denote continued fractions with both unitary numerators and integer coefficients. We will find out that regular continued fractions are enough to represent every real number.

Above, we defined continued fractions as having a finite number of coefficients. We call these *finite continued fractions*. It turns out it also makes sense to talk about fractions that keep going indefinitely, where we have an infinite number of coefficients. Such fractions are called *infinite continued fractions*.

There is a link between continued fractions and the euclidean algorithm, which we will explore by first detouring into an overview of the Gauss map. This will also serve as a prelude to the coming sections, as it our first example of a dynamical system. Define the Gauss map by

$$T(x) = \begin{cases} \left\{ \frac{1}{x} \right\}, & \text{when } 0 < x \leq 1 \\ 0, & \text{when } x = 0. \end{cases} \quad (1)$$

For example $T(\frac{4}{9}) = \{2.25\} = 0.25$. In the section on diophantine approximation we defined notation for the decomposition of a number into integer and fractional part by $\alpha = [\alpha] + \{\alpha\}$. For continued fractions, this is equivalent to $[a_0; a_1, a_2, \dots] = a_0 + [0; a_1, a_2, \dots] =$

$a_0 + 1/[a_1; a_2, \dots]$. Let us try applying the Gauss map over an infinite continued fraction, $\alpha = [0; a_1, a_2, a_3, \dots]$ and see what happens.

$$T(\alpha) = \left\{ a_1 + \frac{1}{a_2 + \frac{1}{a_3 + \frac{1}{\ddots}}} \right\} = [0; a_2, a_3, \dots].$$

Since the leftmost coefficient of α is simply its integer part, the map seems to have shifted the coefficients to the left, removing a_1 in the process. Essentially, we have discovered that we can find the continued fraction coefficients by taking the integer part of the reciprocal of the Gauss map, $a_k = \lfloor 1/T^k(x) \rfloor$. This suggests a method for calculating the continued fraction coefficients of any real number, since the Gauss map is not dependent on α being in the form of a continued fraction. For numbers larger than one, we can remove the integer part and use it as a_0 after the calculations.

So where does the euclidean algorithm come in? Given two natural numbers, say $p \geq q$, the algorithm finds the greatest common divisor (gcd) by iteratively decomposing the p into a maximal product of q and a positive remainder, r . Typically the following scheme is used

$$\begin{aligned} p &= a_0 q + r_0, \\ q &= a_1 r_0 + r_1, \\ &\vdots \\ r_{k-2} &= a_k r_{k-1} + r_k, \\ &\vdots \\ r_{n-3} &= a_{n-1} r_{n-2} + r_{n-1}, \\ r_{n-2} &= a_n r_{n-1} + 0. \end{aligned}$$

It then terminates when $r_n = 0$, giving $\gcd(p, q) = r_{n-1}$. Again, the r_k are the remainders after division. By dividing both sides by the right factor, we can see that $\frac{p}{q} = \left[\frac{p}{q} \right] + \left\{ \frac{p}{q} \right\} = a_0 + \frac{r_0}{q}$. Taking particular note of $\left\{ \frac{p}{q} \right\} = \frac{r_0}{q}$, and checking the next row of the scheme, we find

$$\left\{ \frac{q}{r_0} \right\} = \left\{ \frac{1}{\left\{ \frac{p}{q} \right\}} \right\} = T \left(\left\{ \frac{p}{q} \right\} \right) = \frac{r_1}{r_0}.$$

And in general, $\left\{ \frac{r_{k-2}}{r_{k-1}} \right\} = T^k \left(\left\{ \frac{p}{q} \right\} \right) = \frac{r_k}{r_{k-1}}$. Each step is arrived at from the fractional part of the reciprocal of the last, which is exactly the Gauss map. Naturally, we can as before take the integer part of

the map to get the continued fraction coefficients, putting the intertwinement of continued fractions, the Gauss map, and the euclidean algorithm on full display.

As a side note, let us look at another interesting facet of the euclidean algorithm, namely what happens if the remainders, r_k , are substituted backwards through the scheme. Starting from the second-to-last row of the scheme, $r_{n-3} = a_{n-1}r_{n-2} + r_{n-1}$, we find

$$\begin{aligned} \gcd(p, q) = r_{n-1} &= r_{n-3} - a_{n-1}r_{n-2} \\ &= r_{n-3} - a_{n-1}(r_{n-4} - a_{n-2}r_{n-3}) \\ &= (1 + a_{n-2})r_{n-3} - a_{n-1}r_{n-4} \\ &= (1 + a_{n-2})(r_{n-5} - a_{n-3}r_{n-4}) - a_{n-1}r_{n-4} \\ &= (1 + a_{n-2})r_{n-5} - (a_{n-1} + a_{n-3})r_{n-4} \\ &\vdots \end{aligned}$$

It is hopefully evident from the above calculations that it is possible to always have only two terms consisting of neighbouring remainders r_k , by consistently substituting the r_k with a larger index. This means that by identifying $p = r_{-2}$ and $q = r_{-1}$, we will eventually have something of the form $\gcd(p, q) = mp + nq$, with $m, n \in \mathbb{Z} \setminus \{0\}$. This is known as Bézout's identity, or Bézout's theorem.

Theorem 3.0.1 (Bézout). *Given $a, b \in \mathbb{Z}$, there exists $m, n \in \mathbb{Z}$ such that*

$$\gcd(a, b) = am + bn.$$

With this we are ready to track back, and show that finite continued fractions can be thought of as ordinary fractions.

Theorem 3.0.2. *Finite continued fractions are rational numbers. Additionally, for every rational number there exists exactly two different continued fractions, which are finite.*

Proof. Consider the steps in a slightly rewritten form of the euclidean algorithm for some fraction p/q ,

$$\begin{aligned} \frac{p}{q} &= a_0 + \frac{r_0}{q}, \\ \frac{q}{r_0} &= a_1 + \frac{r_1}{r_0}, \\ \frac{r_0}{r_1} &= a_2 + \frac{r_2}{r_1}, \\ &\vdots \\ \frac{r_{n-2}}{r_{n-1}} &= a_n + 0. \end{aligned}$$

Let us use this to show that the naming of a_k is appropriate in the context of continued fractions, meaning that they are exactly the continued fraction coefficients.

$$\begin{aligned} \frac{p}{q} &= a_0 + \frac{r_0}{q} = a_0 + \frac{1}{\frac{q}{r_0}} \\ &= a_0 + \frac{1}{a_1 + \frac{r_1}{r_0}} = a_0 + \frac{1}{a_1 + \frac{1}{\frac{r_0}{r_1}}} \\ &= a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \frac{r_2}{r_1}}} = \dots \\ &= a_0 + \frac{1}{a_1 + \frac{1}{\dots + \frac{1}{a_n}}} = [a_0; a_1, a_2, \dots, a_n]. \end{aligned}$$

This echoes our conclusions when we discussed the Gauss map earlier. Since the euclidean algorithm can be used on any p, q , we are always able to find a continued fraction representing their quotient. As a consequence, all rational numbers must have a finite continued fraction representation.

As to the existence of at least two continued fractions for each rational, it is enough to note that

$$[a_0; a_1, \dots, a_n] = [a_0; a_1, \dots, a_n - 1, 1],$$

in other words that at the bottommost fraction we can always shift between a_n and $a_n + \frac{1}{1}$ since both are valid when expressing regular continued fractions.

It remains to prove uniqueness when $a_n \neq 1$. We do this by contradiction. Assume there are two continued fractions a and b , such that $[a_0; a_1, \dots, a_n] = [b_0; b_1, \dots, b_m]$ and $a_n \neq 1$ and $b_m \neq 1$. If $n < m$, but $a_k = b_k$ for all $0 \leq k \leq n$, then $a \neq b$ since $[0; b_{n+1}, \dots, b_m] \neq 0$. Thus we can assume that for some $0 \leq k \leq n$, $[a_0; a_1, \dots, a_{k-1}] = [b_0; b_1, \dots, b_{k-1}]$, and $a_k \neq b_k$. Using this assumption, we may rewrite the equality $a = b$ as

$$a_k + \frac{1}{a_{k+1} + \frac{1}{a_{k+2} + \dots}} = b_k + \frac{1}{b_{k+1} + \frac{1}{b_{k+2} + \dots}}.$$

But a_k and b_k are integers, and both right-hand terms are positive and strictly less than one (as long as $b_{k+1} \neq 1$ or $k+1 \neq m$, which is implied by $b_m \neq 1$), so we must have $a_k = b_k$ for the equality to hold. This is contrary to the initial assumption, so for all k , $a_k = b_k$. Then both forms must be equal and have the same number of coefficients. \square

A corollary to the above theorem is that all representations of a fraction share the same continued fraction expansion. The question might then be, which representation does the continued fraction simplify to?

It turns out to be the most simple one, in which the numerator and denominator does not share any common factors.

Theorem 3.0.3. *The numerator and denominator of a finite continued fraction are relatively prime.*

Proof. Consider the fractional form of a finite continued fraction,

$$\alpha = a_0 + \frac{1}{a_1 + \frac{1}{\ddots a_{n-2} + \frac{1}{a_{n-1} + \frac{1}{a_n}}}}.$$

By rewriting the bottommost addition we get $\frac{a_{n-1}a_n+1}{a_n}$, the numerator and denominator of which are easily seen to be relatively prime. If we substitute in $s = a_{n-1}a_n + 1$ and $t = a_n$, then doing the next addition up in the continued fraction, we get $a_{n-2} + \frac{t}{s} = \frac{a_{n-2}s+t}{s}$. This is again a coprime fraction, since s and t do not have common factors. We can repeat the same argument until we are left with an ordinary fraction, which demonstrably have numerator and denominator relatively prime. \square

The continued fraction coefficients form a sequence, so we would like to be able to make statements about the number as a whole from only taking the first few coefficients. This will in turn allow us to work inductively on both finite and infinite continued fractions.

Definition 3.0.4. The k -convergent of a regular continued fraction $\alpha = [a_0; a_1, a_2, \dots]$ is the rational corresponding to $[a_0; a_1, a_2, \dots, a_k]$.

Finite continued fractions are rational, and it makes sense to define p_k and q_k , such that their quotient gives the corresponding convergent, $\frac{p_k}{q_k} = [a_0; a_1, a_2, \dots, a_k]$. We call p_k partial numerators and q_k partial denominators.

The following two theorems are very useful when working with k -convergents.

Theorem 3.0.5. *For all $k \geq 2$,*

$$p_k = a_k p_{k-1} + p_{k-2},$$

$$q_k = a_k q_{k-1} + q_{k-2}.$$

Proof. Using induction, for a continued fraction $\alpha = [a_0; a_1, a_2, \dots]$, we have

$$\frac{p_0}{q_0} = \frac{a_0}{1},$$

$$\frac{p_1}{q_1} = a_0 + \frac{1}{a_1} = \frac{a_0 a_1 + 1}{a_1},$$

$$\frac{p_2}{q_2} = a_0 + \frac{1}{a_1 + \frac{1}{a_2}} = \frac{a_0 a_1 a_2 + a_0 + a_2}{a_1 a_2 + 1}.$$

We confirm that

$$a_2 p_1 + p_0 = a_2 a_0 a_1 + a_2 + a_0 = p_2,$$

and

$$a_2 q_1 + q_0 = a_2 a_1 + 1 = q_2.$$

For the general case, introduce $\frac{\hat{p}_k}{\hat{q}_k} = [a_1; a_2, \dots, a_k]$. Then

$$\begin{aligned} \frac{p_k}{q_k} &= a_0 + \frac{1}{\frac{\hat{p}_k}{\hat{q}_k}} \\ &= a_0 + \frac{1}{(a_k \hat{p}_{k-1} + \hat{p}_{k-2}) / (a_k \hat{q}_{k-1} + \hat{q}_{k-2})} \\ &= \frac{a_k (a_0 \hat{p}_{k-1} + \hat{q}_{k-1}) + a_0 \hat{p}_{k-2} + \hat{q}_{k-2}}{a_k \hat{p}_{k-1} + \hat{p}_{k-2}}. \end{aligned}$$

We will need to show that $p_k = a_0 \hat{p}_k + \hat{q}_k$ and $q_k = \hat{p}_k$. Again,

$$\frac{p_k}{q_k} = a_0 + \frac{1}{\frac{\hat{p}_k}{\hat{q}_k}} = \frac{a_0 \hat{p}_k + \hat{q}_k}{\hat{p}_k}.$$

So $p_k = \lambda(a_0 \hat{p}_k + \hat{q}_k)$ and $q_k = \lambda \hat{p}_k$, and since p_k and q_k are relatively prime, $\lambda = 1$. Finally,

$$\begin{aligned} \frac{p_k}{q_k} &= \frac{a_k (a_0 \hat{p}_{k-1} + \hat{q}_{k-1}) + a_0 \hat{p}_{k-2} + \hat{q}_{k-2}}{a_k \hat{p}_{k-1} + \hat{p}_{k-2}} \\ &= \frac{\lambda (a_k p_{k-1} + p_{k-2})}{\lambda (a_k q_{k-1} + q_{k-2})} = \frac{a_k p_{k-1} + p_{k-2}}{a_k q_{k-1} + q_{k-2}}. \end{aligned}$$

□

Theorem 3.0.6. For $k \geq 1$,

$$\frac{p_{k-1}}{q_{k-1}} - \frac{p_k}{q_k} = \frac{(-1)^k}{q_{k-1} q_k}.$$

Proof. First note the equivalent form,

$$p_{k-1} q_k - q_{k-1} p_k = (-1)^k.$$

We prove that it holds with induction. First show it holds for $k = 1$: $p_0 = a_0$, $p_1 = a_0 a_1 + 1$, $q_0 = 1$, $q_1 = a_1$. So

$$a_0 a_1 - 1 \cdot (a_0 a_1 + 1) = -1.$$

Now assume $p_{k-2} q_{k-1} - q_{k-2} p_{k-1} = (-1)^{k-1}$. Then

$$\begin{aligned} p_{k-1} q_k - q_{k-1} p_k &= p_{k-1} (a_k q_{k-1} + q_{k-2}) \\ &\quad - q_{k-1} (a_k p_{k-1} + p_{k-2}) \\ &= p_{k-1} q_{k-2} - q_{k-1} p_{k-2} \\ &\quad + a_k (p_{k-1} q_{k-1} - q_{k-1} p_{k-1}) \\ &= -(-1)^{k-1} + 0 = (-1)^k. \end{aligned}$$

□

Now we show that the $\frac{p_k}{q_k}$ are well behaved in terms of convergence.

Theorem 3.0.7. *The sequence of an arbitrary infinite continued fraction's k -convergents converges to a real number α .*

For the proof we will utilise the following lemma.

Lemma 3.0.8. *The partial numerators and denominators of a (non-zero) continued fraction are bounded by the Fibonacci numbers according to*

$$|p_k| \geq F_k \text{ and } q_k \geq F_{k+1}.$$

Here we use the convention $F_0 = 0$, and $F_1 = 1$.

Proof. We can prove this by induction. First we confirm for $k = 0$ and $k = 1$ that $p_0 = a_0$, $p_1 = a_0 a_1 + 1$, $q_0 = 1$, $q_1 = a_1$.

$$|p_0| = |a_0| \geq 0 = F_0, \quad |p_1| = |a_0 a_1 + 1| \geq 1 = F_1.$$

$$q_0 = 1 \geq 1 = F_1, \quad q_1 = a_1 \geq 1 = F_2.$$

Now assume it holds for $k-2$ and $k-1$. Because the integer part of a continued fraction is given by a_0 , and the fractional part is positive, all partial numerators p_k have the same sign as a_0 . Thus, for $k \geq 2$,

$$\begin{aligned} |p_k| &= |a_k p_{k-1} + p_{k-2}| = a_k |p_{k-1}| + |p_{k-2}| \\ &\geq a_k F_{k-1} + F_{k-2} \geq F_{k-1} + F_{k-2} = F_k. \end{aligned}$$

Similarly, since all partial denominators q_k are positive,

$$q_k = a_k q_{k-1} + q_{k-2} \geq a_k F_k + F_{k-1} \geq F_{k+1}.$$

□

Proof of 3.0.7. Let $[a_0; a_1, \dots]$ be an infinite regular continued fraction. Using Theorem 3.0.6 together with the above lemma, we get

$$\left| \frac{p_{k-1}}{q_{k-1}} - \frac{p_k}{q_k} \right| = \frac{1}{q_{k-1} q_k} \leq \frac{1}{F_k F_{k+1}}.$$

Since for $k \geq 5$, we have $F_{k+1} > F_k \geq k$, then

$$\sum_{k=1}^{\infty} \left| \frac{p_{k-1}}{q_{k-1}} - \frac{p_k}{q_k} \right| \leq \sum_{k=1}^{\infty} \frac{1}{F_k F_{k+1}} \leq C + \sum_{k=1}^{\infty} \frac{1}{k^2} < \infty,$$

where C is some finite constant. This means that $\frac{p_k}{q_k}$ is a Cauchy sequence, and therefore converges to a real number. □

The converse also holds.

Theorem 3.0.9. *For every irrational number there is an unique regular continued fraction whose k -convergents converges to it.*

Proof. Using the euclidean algorithm to generate a continued fraction for α , we get an infinite number of coefficients. We also have a nonregular description $\alpha = [a_0; \dots, a_{k-1}, r_k]$, where r_k is a remainder > 1 , such that $\lfloor r_k \rfloor = a_k$. By Theorem 3.0.5, we find $\frac{p_k}{q_k} = \frac{p_{k-1}a_k + p_{k-2}}{a_k q_{k-1} + q_{k-2}}$ and $\alpha = \frac{r_k p_{k-1} + p_{k-2}}{r_k q_{k-1} + q_{k-2}}$. Computing the difference gives us

$$\begin{aligned} \left| \alpha - \frac{p_k}{q_k} \right| &= \left| \frac{r_k p_{k-1} + p_{k-2}}{r_k q_{k-1} + q_{k-2}} - \frac{a_k p_{k-1} + p_{k-2}}{a_k q_{k-1} + q_{k-2}} \right| \\ &= \left| \frac{(p_{k-1} q_{k-2} - p_{k-2} q_{k-1})(r_k - a_k)}{(r_k q_{k-1} + q_{k-2})(a_k q_{k-1} + q_{k-2})} \right| \\ &= \left| \frac{(-1)^k (r_k - a_k)}{(r_k q_{k-1} + q_{k-2})(a_k q_{k-1} + q_{k-2})} \right| \\ &< \frac{1}{q_k^2} \leq \frac{1}{F_k^2}. \end{aligned}$$

Note that we in the second to last step used that

$$a_k < r_k \Rightarrow r_k q_{k-1} + q_{k-2} > a_k q_{k-1} + q_{k-2} = q_k.$$

This shows that the difference goes to 0 as k goes to infinity, so $\frac{p_k}{q_k}$ converges to α .

We can show uniqueness in a way similar to what we did for the finite case. Assume there is another continued fraction $\beta = [b_0; b_1, \dots]$ converging to α . Then there is a first k for which $a_k \neq b_k$. But convergence means $[a_0; a_1, \dots] = [b_0; b_1, \dots]$, implying

$$a_k + \frac{1}{a_{k+1} + \frac{1}{\ddots}} = b_k + \frac{1}{b_{k+1} + \frac{1}{\ddots}}.$$

This cannot be true since the fractional terms are smaller than 1. So $a_k = b_k$, going against our assumption that α and β are different. \square

With the help of k -convergents, we are able to make statements about approximability.

Theorem 3.0.10. *The coefficients of an infinite continued fraction is eventually periodic iff it converges to a quadratic irrational.*

The theorem can be naturally divided into two halves, the first being that eventually periodic continued fractions are quadratic irrationals.

Proof of the first half of 3.0.10. Introduce $\alpha = [\overline{a_0; a_1, a_2, \dots, a_n}]$, which is the periodic continued fraction expansion of a real number, as well as $\beta = [b_0; b_1, b_2, \dots, b_m, \alpha]$, an eventually periodic continued fraction containing α . Note that β is in regular form, despite having an irrational coefficient, since α can be substituted for its coefficients. Making use of Theorem 3.0.5 we find

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

(This assumes that the theorem holds for non-integer terminal coefficients.) With some rewriting, this becomes

$$q_n^\alpha \alpha^2 + (q_{n-1}^\alpha - p_n^\alpha) \alpha - p_{n-1}^\alpha = 0,$$

revealing that α is the root of a quadratic equation. In a similar fashion,

$$\beta = \frac{\alpha p_m^\beta + p_{m-1}^\beta}{\alpha q_m^\beta + q_{m-1}^\beta}.$$

From this we are able to express α in terms of β ,

$$\alpha = -\frac{p_{m-1}^\beta - \beta q_{m-1}^\beta}{p_m^\beta - \beta q_m^\beta},$$

which we can then substitute back into the quadratic equation for α , yielding

$$\begin{aligned} q_n^\alpha (p_{m-1}^\beta - \beta q_{m-1}^\beta)^2 \\ + (q_{n-1}^\alpha - p_n^\alpha) (p_{m-1}^\beta - \beta q_{m-1}^\beta) (p_m^\beta - \beta q_m^\beta) \\ + p_{n-1}^\alpha (p_m^\beta - \beta q_m^\beta)^2 = 0 \end{aligned}$$

The partial numerators and denominators are integers, and we can see that β is at most squared, so β is the root of a quadratic equation. \square

The second part of the theorem, which can be attributed to Lagrange, is that quadratic irrationals have periodic continued fraction expansions, and to prove it we will require some intermediate results. The first thing we shall need is the notion of *residue*. As we are working with real roots of quadratic equations with integer coefficients, we can assume that α is a positive irrational number greater than 1. We will write $\alpha = a_0 + \frac{1}{\alpha_1}$, where $\alpha_1 = [a_1; a_2, a_3, \dots]$ is the first residue of α . In general, we may define the k -th residue of α as $\alpha_k = a_k + \frac{1}{\alpha_{k+1}}$. This leads us to the following lemma.

Lemma 3.0.11. *Let α be a positive irrational number greater than 1, solving the equation*

$$rx^2 + sx + t = 0,$$

for some $r, s, t \in \mathbb{Z}$. Then the residues α_k solve similar quadratic equations,

$$r_k x^2 + s_k x + t_k = 0,$$

such that $t_{k+1} = r_k$, and the discriminants $s_k^2 - 4r_k t_k = s^2 - 4rt$ are independent of k .

Proof. The idea is to find t_1, s_1, r_1 solving the equation

$$t_1(x - a_0)^2 + s_1(x - a_0) + r_1 = 0,$$

since then the residue $\alpha_1 = \frac{1}{\alpha - a_0}$ satisfies

$$r_1 \alpha_1^2 + s_1 \alpha_1 + t_1 = 0,$$

which can be seen by multiplying in $(x - a_0)^{-4}$. One way to find the coefficients is by equating the the above equation to the one in the lemma statement, giving

$$t_1(x - a_0)^2 + s_1(x - a_0) + r_1 = rx^2 + sx + t.$$

After some shuffling around of the above, we can eventually identify $t_1 = r$, $s_1 = s + 2ra_0$ and $r_1 = ra_0^2 + sa_0 + t$. With the residue coefficients known, we are also able to confirm that $s_1^2 - 4r_1t_1 = s^2 - 4rt$. With the equation

$$r_1x^2 + s_1x + t_1 = 0,$$

having the same form as the one in the lemma statement, we can repeat the same procedure any number of times to find $t_{k+1} = r_k$, $s_{k+1} = s_k + 2r_k a_k$ and $r_{k+1} = r_k a_k^2 + s_k a_k + t_k$. \square

We will also need a certain property of the r_k .

Lemma 3.0.12. *The sequence of r_k , defined in the previous lemma, switches signs an infinite number of times.*

Proof. Assume instead that the sequence switches signs a finite number of times. Then there exists an n for which all $r_k, k \geq n$ have the same sign. Since $t_{k+1} = r_k$, all t_k after that point must have the same sign as well. Now recall $s_{k+1} = s_k + 2r_k a_k$. For $k \geq n$ we can substitute $s_k = s_{k-1} + 2r_{k-1} a_{k-1}$ into the equation until we find

$$s_{k+1} = s_n + \sum_{m=n}^k 2r_m a_m.$$

Since all a_k are positive integers, and r_k for $k \geq n$ are integers with the same sign, then for k sufficiently big, the sum $\sum_{m=n}^k 2r_m a_m$ must eventually have an absolute value greater than s_n . After that point all s_k have the same sign as r_k . So for $k \gg n$, we must have that r_k, s_k and t_k all have the same sign. However since

$$r_k x^2 + s_k x + t_k = 0,$$

has no positive real solutions when all coefficients have the same sign, then contrary to construction $\alpha_k \in \mathbb{R}$ cannot be one of the roots. We conclude that a finite number of sign changes leads to a contradiction. \square

We now have the tools to prove that quadratic irrationals have periodic continued fraction expansions.

Proof of the second half of 3.0.10 (Lagrange). We have proved in the first lemma that the discriminant is independent of k ,

$$s_k^2 - 4r_k t_k = s^2 - 4rt.$$

From the second lemma, and the fact $t_{k+1} = r_k$, we know that there is an infinite number of times that $4r_k t_k$ is negative. This means $s^2 - 4rt$ is described as the sum of two positive integers an infinite number of times. But there is only a finite number of ways to describe a natural number in such a way, so at some point there must be a repetition of the coefficients in the residue equations. Since r_k, s_k, t_k uniquely define $r_{k+1}, s_{k+1}, t_{k+1}$ and α_{k+1} , then after that repetition the coefficients must repeat periodically. \square

By now we have done quite a bit of work showing different aspects of the continued fractions, where the main takeaway should be that fractions on this form are well behaved and an efficient representation of diophantine approximation. With that in mind, let us revisit some of the concepts from our discussion on approximation in the first section.

Theorem 3.0.13. *All k -convergents for an irrational number α are best approximation of α .*

For the proof we shall use a property inherent to the way the convergents of α converge, namely that they do so by overshooting.

Lemma 3.0.14. *Even-numbered k -convergents approach α from below, that is*

$$\frac{p_{2k}}{q_{2k}} < \frac{p_{2k+2}}{q_{2k+2}} < \alpha.$$

Odd-numbered k -convergents approach α from above, that is

$$\frac{p_{2k-1}}{q_{2k-1}} > \frac{p_{2k+1}}{q_{2k+1}} > \alpha.$$

Proof. From Theorem 3.0.6 we know

$$\frac{p_{k-1}}{q_{k-1}} - \frac{p_k}{q_k} = \frac{(-1)^k}{q_{k-1}q_k}.$$

As q_k is positive for all k , we see that $\frac{p_0}{q_0} < \frac{p_1}{q_1}, \frac{p_1}{q_1} > \frac{p_2}{q_2}, \dots, \frac{p_{2k}}{q_{2k}} < \frac{p_{2k+1}}{q_{2k+1}}$.

Using theorems 3.0.5 and 3.0.6 we can also relate $\frac{p_{k+2}}{q_{k+2}}$ to $\frac{p_k}{q_k}$,

$$\begin{aligned}\frac{p_{k+2}}{q_{k+2}} &= \frac{p_{k+1}}{q_{k+1}} - \frac{(-1)^k}{q_{k+1}q_{k+2}} \\ &= \frac{p_k}{q_k} + \frac{(-1)^k}{q_k q_{k+1}} - \frac{(-1)^k}{q_{k+1}q_{k+2}} \\ &= \frac{p_k}{q_k} + (-1)^k \frac{q_{k+2} - q_k}{q_k q_{k+1} q_{k+2}} \\ &= \frac{p_k}{q_k} + (-1)^k \frac{\alpha_{k+2} q_{k+1} + q_k - q_k}{q_k q_{k+1} q_{k+2}} = \frac{p_k}{q_k} + (-1)^k \frac{\alpha_{k+2}}{q_k q_{k+2}}.\end{aligned}$$

Again, q_k is always positive, as is α_k , so the sequence of $\frac{p_{2k}}{q_{2k}}$ is strictly increasing, and $\frac{p_{2k+1}}{q_{2k+1}}$ is strictly decreasing. With these two observations we find

$$\frac{p_{2k}}{q_{2k}} < \frac{p_{2k+2l}}{q_{2k+2l}} < \frac{p_{2k+2l+1}}{q_{2k+2l+1}} < \frac{p_{2l+1}}{q_{2l+1}},$$

for all k and l . The result follows since we know that the k -convergents converge to α . \square

Proof of Theorem 3.0.13. Suppose we have some approximation of α , such that $\left| \alpha - \frac{p}{q} \right| < \left| \alpha - \frac{p_k}{q_k} \right|$ for some particular k . Then $\frac{p_k}{q_k}$ is only a best approximation of α if $q \geq q_k$, which is what we will show. For the sake of convenience, let us assume that k is odd, although a similar argument can be made if k is even. Then according to the above lemma,

$$\frac{p_{k-1}}{q_{k-1}} < \frac{p}{q} < \frac{p_k}{q_k},$$

which we can rewrite as

$$0 < \frac{p}{q} - \frac{p_{k-1}}{q_{k-1}} < \frac{p_k}{q_k} - \frac{p_{k-1}}{q_{k-1}},$$

and again

$$0 < \frac{p q_{k-1} - p_{k-1} q}{q q_{k-1}} < \frac{p_k q_{k-1} - p_{k-1} q_k}{q_k q_{k-1}}.$$

Using Theorem 3.0.6 on the rightmost expression and that the numerators are positive integers, we finally get

$$\frac{1}{q q_{k-1}} < \frac{1}{q_k q_{k-1}},$$

or in other words, $q > q_k$. \square

So the convergents of α are best approximations. However we earlier defined best approximations as fractions $\frac{p^*}{q^*}$ for which

$$\left| \alpha - \frac{p^*}{q^*} \right| < \left| \alpha - \frac{p}{q} \right|$$

for all $p < p^*$ and $q < q^*$. These are more specifically best approximations of the *first kind*. Approximations of this kind also include rationals not in the convergents of α . For example, let us return to π , for which the first few regular continued fraction coefficients are

$$\pi = [3; 7, 15, 1, 292, 1, 1, 1, 2, 1, 3, 1, 14, 2, 1, 1, 2, 2, 2, 2, 1, 84, \dots],$$

where the approximation $\pi \approx 3$ is given by only taking the first coefficient, and $\pi \approx \frac{22}{7} = 3 + \frac{1}{7}$ is given by the first two coefficients, $[3; 7]$. Since $\left| \pi - \frac{13}{4} \right| \approx 0.11 < 0.14$, clearly there is a best approximation between 3 and $\frac{22}{7}$. Then what is special about the convergents? This is where the difference between the forms $\left| \alpha - \frac{p}{q} \right|$ and $\|q\alpha\|$ comes in.

Definition 3.0.15. A best approximation of the *second kind* of a real number α is a positive integer q^* such that for all $q < q^*$,

$$\|q^*\alpha\| < \|q\alpha\|.$$

By this definition of best approximations, the approximation $\pi \approx \frac{13}{4}$ gives $\|4\pi\| \approx 0.43$. This is not better than $\frac{3}{1}$, for which $\|1\pi\| \approx 0.14$. Meanwhile $\pi \approx \frac{22}{7}$ is still considered a best approximation with $\|7\pi\| \approx 0.0089$. In fact, the convergents of α are exactly the best approximations of the second kind, and in this sense the form $\|q\alpha\|$ is slightly stronger as all approximations of the second kind are also that of the first kind.

The next theorem seems like an appropriate ending to this chapter, as it makes use of most things we have learned up to this point, and will have some bearing on our understanding of our discussions later.

Theorem 3.0.16. *An irrational number, α , is badly approximable iff the coefficients of its continued fraction expansion are bounded. Since periodicity implies boundedness, this can be considered as an extension of Theorem 3.0.10.*

Proof. Let us begin by relating $\|q\alpha\|$ to the convergents of α , using the nonstandard continued fraction form,

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

where α_{n+1} is a real number. Then

$$\begin{aligned} \left| \alpha - \frac{p_n}{q_n} \right| &= \left| \frac{\alpha_{n+1}p_n - p_{n-1} - \frac{p_n}{q_n}}{\alpha_{n+1}q_n - q_{n-1}} \right| \\ &= \left| \frac{\alpha_{n+1}p_nq_n + p_{n-1}q_n - \alpha_{n+1}p_nq_n - p_nq_{n-1}}{q_n(\alpha_{n+1}q_n + q_{n-1})} \right| \\ &= \left| \frac{p_{n-1}q_n - p_nq_{n-1}}{q_n(\alpha_{n+1}q_n + q_{n-1})} \right| = \frac{1}{q_n(\alpha_{n+1}q_n + q_{n-1})}. \end{aligned}$$

Now we multiply by q_n , and get

$$\|q_n\alpha\| = \frac{1}{\alpha_{n+1}q_n + q_{n-1}} < \frac{1}{a_{n+1}q_n}.$$

For the inequality, recall that $\alpha_{n+1} \geq [\alpha_{n+1}] = a_{n+1}$. Using the definition of badly approximable numbers, then α is badly approximable if and only if there exists a $c > 0$, such that we for all $q_n \in \mathbb{Z}$ have

$$\frac{c}{q_n} < \|q_n \alpha\|.$$

Now $\frac{c}{q_n} < \frac{1}{a_{n+1} q_n}$, hence a_{n+1} is bounded by $1/c$.

That boundedness is a sufficient condition should become clear if we consider best approximations. Combining Theorem 3.0.13 and Lemma 3.0.14 we surmise that α lies closer to $\frac{p_{n+1}}{q_{n+1}}$ than $\frac{p_n}{q_n}$, and in particular that their midpoint is smaller than α , meaning

$$\frac{p_{n+1}/q_{n+1} + p_n/q_n}{2} < \alpha.$$

As the left-hand side is positive, we are by subtracting $\frac{p_n}{q_n}$ and taking absolutes able to simplify the inequality to

$$\left| \frac{p_{n+1}/q_{n+1} - p_n/q_n}{2} \right| < \left| \alpha - \frac{p_n}{q_n} \right|.$$

Next we note that $|p_{n+1}/q_{n+1} - p_n/q_n| = \frac{1}{q_n q_{n+1}}$, therefore, after multiplication by q_n we can rewrite the inequality,

$$\frac{1}{2q_{n+1}} < \|q_n \alpha\|.$$

In particular, $\frac{p_n}{q_n}$ being a best approximation means for all $q < q_n$,

$$\frac{1}{2(a_{n+1} + 1)q_n} < \frac{1}{2(a_{n+1}q_n + q_{n-1})} = \frac{1}{2q_{n+1}} < \|q_n \alpha\| < \|q\alpha\|.$$

Since the coefficients are bounded from above, there exists an integer $M > \sup_n a_n$, and we may therefore choose $c = \frac{1}{2(M+1)}$, and α is badly approximable. \square

3.0.0.1 Additional notes

The euclidean algorithm, and therefore the calculation of continued fractions, may be restated in matrix form. Consider again the steps in the algorithm for some irreducible fraction p/q ,

$$\begin{aligned} p &= qa_0 + r_0, \\ q &= r_0a_1 + r_1, \\ r_0 &= r_1a_2 + r_2, \\ &\vdots \\ r_{n-2} &= r_{n-1}a_n + 0. \end{aligned}$$

In general,

$$r_k = r_{k+1} a_{k+2} + r_{k+2}.$$

We can get a system of equations by stating the tautology

$$r_{k+1} = r_{k+1} + 0 \cdot r_{k+2}.$$

Then, in matrix form we have

$$\begin{pmatrix} r_k \\ r_{k+1} \end{pmatrix} = \begin{pmatrix} a_{k+2} & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} r_{k+1} \\ r_{k+2} \end{pmatrix}.$$

Starting from $\begin{pmatrix} p \\ q \end{pmatrix}$, we can then find a matrix expression of the euclidean algorithm, by iteratively rewriting the vector

$$\begin{pmatrix} p \\ q \end{pmatrix} = \begin{pmatrix} a_0 & 1 \\ 1 & 0 \end{pmatrix} \cdots \begin{pmatrix} a_n & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

This form is particularly interesting considering its relation to Möbius transformations, which we will cover in Section 5.

An example of a continued fraction representation is that of the golden ratio, $\phi = [1; 1, 1, 1, \dots]$. It was mentioned that ϕ is the worst approximable number, which can be attributed to the fact that the continued fraction coefficients are the smallest possible, and in general that large continued fraction coefficients mean the corresponding k -convergent is a large improvement over the previous convergent in terms of approximating α .

While there does not seem to be a tractable pattern to the continued fraction coefficients of π , there are several more well behaved continued fractions if the restriction of regularity is dropped. One such fraction is

$$\pi = \frac{4}{1 + \frac{1^2}{3 + \frac{2^2}{5 + \frac{3^2}{\ddots}}}}$$

where the coefficient places are filled by the odd numbers, and the numerators are given by the squares of the natural numbers.

The Gauss map is ergodic, and it is possible to calculate the distribution of a_k . Recall that the coefficients of a real number α could be calculated by repeated application of the Gauss map, i.e. $a_n = \lfloor T^n(\alpha) \rfloor$. From this we are able to define the measure $m_n(x)$ over the set of real numbers $\alpha \in [0, 1]$ where $T^n(\alpha) < x$. Then we introduce the following theorem

Theorem 3.0.17 (Gauss-Kuzmin). For $0 \leq x \leq 1$,

$$\lim_{n \rightarrow \infty} m_n(x) = \frac{\ln(1+x)}{\ln 2}.$$

We may from this calculate the probability that $a_n = k$, denoted $P(k)$ as

$$P(k) = \lim_{n \rightarrow \infty} m_n\left(\frac{1}{k}\right) - m_n\left(\frac{1}{k+1}\right) = \frac{1}{\ln 2} \ln\left(1 + \frac{1}{k(k+2)}\right).$$

With the use of ergodic theory it is then possible to prove that the coefficients for almost all real numbers have the same geometric mean,

$$K_0 = \lim_{n \rightarrow \infty} (a_0 a_1 \cdots a_n)^{1/n}.$$

called Khinchin's constant.

SYMBOLIC SEQUENCES

In the following sections we will examine two different systems relating back to continued fractions. The first has its roots in symbolic sequences.

Definition 4.0.1. A *binary symbolic sequence* is a sequence that can be represented as a sequence of zeroes and ones,

$$\sigma = \sigma_1 \sigma_2 \dots$$

where $\sigma_k \in \{0, 1\}$.

We call contiguous subsequences of σ *words*. For example, given that $\sigma = 0010100\dots$, we say that 101 is a word of σ with length 3.

A symbolic sequence is said to be *periodic* if for some $\delta > 0$, and for all n , we have $\sigma_n = \sigma_{n+\delta}$. The period of the sequence is then the smallest such δ that satisfies the equality.

We are particularly interested in a specific subset of the binary symbolic sequences which are called the Sturmian sequences. They may be characterised by their complexity.

Definition 4.0.2. The complexity of a symbolic sequence σ is given by the complexity function $p_\sigma(n)$, which is the number of different words of length n in σ . By convention we use $p_\sigma(0) = 1$, indicating the empty word. As an example, take the periodic sequence $\sigma = 001001001\dots$. Then we can enumerate all subwords of length 2: 00, 01, and 10. Therefore, $p_\sigma(2) = 3$.

The complexity of σ is loosely related to encoding (or substituting) the sequence using only words of length n . Again using $\sigma = 001001001\dots$, we can use the encoding $a = 00$, $b = 10$, and $c = 01$, giving $\sigma = abcabcabc\dots$.

Interestingly, if we instead use words of length 3, we can encode the sequence by $a = 001$ due to periodicity. This gives $\sigma = aaa\dots$, even though $p_\sigma(3) = 3$.

Theorem 4.0.3. A sequence is eventually periodic if for some n ,

$$p_\sigma(n+1) = p_\sigma(n).$$

By eventually periodic, we mean that for some $k \geq 1$, $\hat{\sigma} = \sigma_k \sigma_{k+1} \sigma_{k+2} \dots$ is periodic. Additionally, iff σ is eventually periodic, then p_σ is bounded.

Proof. Assume that n is the smallest number such that $p_\sigma(n) = p_\sigma(n+1) = m$ for some sequence σ . Then each word in σ of length $n+1$

is a word of length n , with only one of either a 0 or 1 appended. For example take some sequence starting with $\sigma = 001\dots$, where $p_\sigma(2) = p_\sigma(3) = 3$. Since 001 is 00 with a 1 appended, 000 can never occur, otherwise the complexity would have increased. As a side note, the sequence is still not uniquely defined as both 001001001... and 001010101... fulfill the given complexity requirements. Now consider $p_\sigma(n+2)$. We know that there are exactly m words of length $n+1$. However, we also know that any m consecutive symbols unambiguously spell out the next symbol. So by taking all but the first symbol of a word of length $n+1$, we can conclude that for each such word of length $n+1$, there is exactly one word of length $n+2$. We can inductively repeat the same argument, and so we find that $p_\sigma(n+k) = m$ for all $k \geq 0$. As we move along the sequence, we will eventually return to a word of length n occurring earlier due to the fact that there are only a finite number of words of any particular length. Since n symbols uniquely implies the following symbol, the sequence, from that point on, will periodically repeat.

Conversely, if σ is eventually periodic, then $p_\sigma(n)$ is bounded, as any sequence with period n , by a similar argument to above, has a finite number of words, and any non-periodicity at the start of the sequence is also only able to contribute a finite number of words. \square

As a corollary, a sequence is periodic if for some n , $p_\sigma(n) = n$, implied by the fact that complexity functions are monotonously increasing. Expressed in another way, for the smallest n satisfying $p_\sigma(n) = n$, we have $p_\sigma(n-1) = n$. This is since if $p_\sigma(n-1)$ was smaller, n could not be the first time for which $p_\sigma(n) = n$. If it was larger, $p_\sigma(n)$ could not be equal to n . Boundedness gives the converse, namely if σ is periodic, then $\exists n, p_\sigma(n) = n$. From this, one could imagine there might be sequences where the complexity is always "just out of reach" of periodicity, which is essentially the definition of Sturmian sequences.

Definition 4.0.4. A Sturmian sequence σ is a sequence where for all n , $p_\sigma(n) = n+1$. In other words, they have the smallest possible complexity for each n , without being periodic. Since by definition $p_\sigma(1) = 2$, all Sturmian sequences are binary.

It is not obvious that Sturmian sequences actually exist, but they do in fact turn up in many different contexts. Let us examine one way of generating Sturmian sequences, which is closely related to Diophantine approximations.

Definition 4.0.5. Rotation sequences occur in the circle shift map. Consider such a map, $S_\alpha : S^1 \rightarrow S^1, x \mapsto \{x + \alpha\}$ (recall that $\{x\}$ is used to denote the fractional part of x). The set S^1 should be considered to be the interval $[0, 1)$, isomorphic to the unit circle in the plane. Now divide S^1 into two intervals $I_0 = [0, 1 - \alpha)$ and $I_1 = [1 - \alpha, 1)$, see Figure 1, and build a sequence by iterating $S_\alpha(x)$. When $x \in I_0$, we

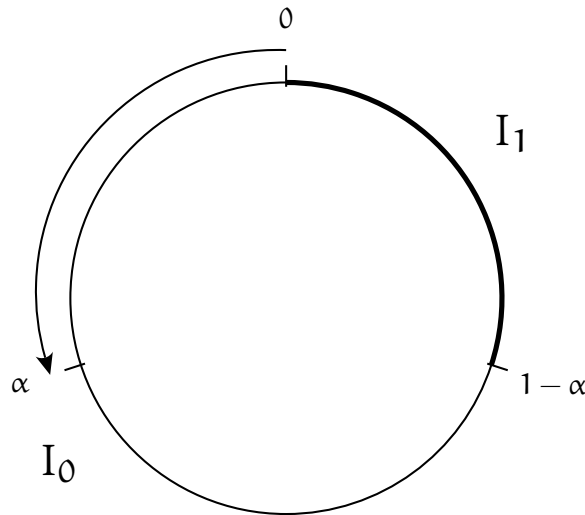


Figure 1: Rotation sequences can be modeled as movements on a circle.

append a 0 to our sequence, and when $x \in I_1$ we append a 1. This then forms the rotation sequence for α .

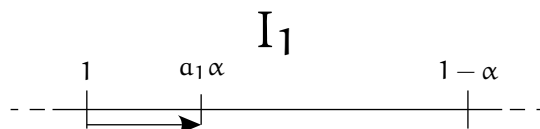
Incidentally, these sequences may also be expressed as

$$\sigma_n = \lfloor (n + 1)\alpha \rfloor - \lfloor n\alpha \rfloor.$$

In general we would leave the starting point as a variable, i.e. $x_0 = \beta$, but for our purposes we will assume $x_0 = 0$. Also, for convenience we will from now on only work with irrational α , since rational values repeat periodically, and are not as interesting as the irrational case.

As alluded to, these sequences are a kind of diophantine approximation. Should we find a $q\alpha$, $q \in \mathbb{N}$, such that $\{q\alpha\} \approx 0$ (or also $\{q\alpha\} \approx 1$), then $p \approx q\alpha$ or $\frac{p}{q} \approx \alpha$ for some p . The circle shift map amounts to a multiplication modulo 1, so $S_\alpha^q(0) = \{q\alpha\}$. With the intent of using this common ground to fully connect diophantine approximation and rotation sequences, we might come up with the following scheme, in which we iteratively find moments in the rotation sequence that get closer and closer to 0 (or equivalently, 1).

We are not interested in the first step, $x_0 = 0$, since this in diophantine approximation corresponds to $q = 0$, so we begin with $x_1 = \alpha$. Let us find the first time that we are closer to 1 than α . We cannot completely step over I_1 using only rotations of length α , so this event coincides with the first time we enter I_1 . Let us write the rotation sequence up to that point as $C_1 = 0^{a_1-1}1$, where a_1 is the number of rotations to get to there. The position on the circle can then be written as $a_1\alpha$, which in turn can be thought of as a rotation in itself.



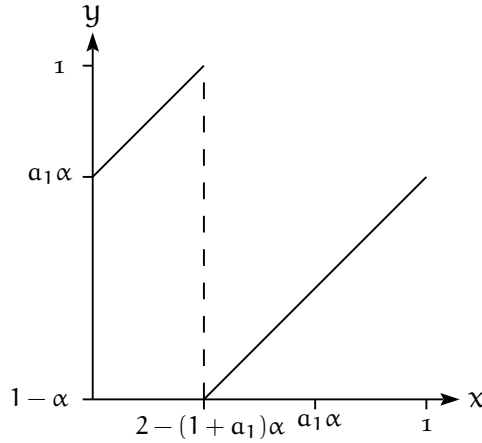


Figure 2: The first return map of I_1 under the circle shift map, $S_\alpha(x)$.

Now we ask, when is the first time in the sequence that we are closer to 1 than $a_1\alpha$? After some thinking we might deduce that it will be at a point in the interval $[0, 1 - a_1\alpha)$, the left half of the interval $(a_1\alpha, 1 - a_1\alpha)$ of all values closer to 1 (or 0) than $a_1\alpha$. Consider in particular where we are after another a_1 steps, which intuitively corresponds to a backwards rotation of length $1 - a_1\alpha$. Either we wind up back in I_1 , or we fall into I_0 , which may be visualised by the first return map of I_1 under S_α , seen in Figure 2. The first return map takes each point in I_1 to the point of its first return to I_1 , subject to iteration under S_α .

We can see two distinct cases. In the case for values in $I'_1 = [1 - \alpha, 2 - (1 + a_1)\alpha)$, since we do not make a complete lap around the circle in a_1 steps, we need one more step to come back into I_1 compared to the other case. However, if we put I'_1 into the shift map, $S_\alpha(I'_1) = [0, 1 - a_1\alpha)$, we note that we in fact get closer than $a_1\alpha$ in only 1 iteration. The points in $I'_0 = [2 - (1 + a_1)\alpha, 1)$ on the other hand will return to I_1 in exactly a_1 steps, and must therefore have landed further away from 1 when they come back to I_1 . We must also not have landed closer at any point on our walk through I_0 , since that is covered by the points in I'_1 . So we can be content that the points in I'_0 do not land closer to 0 during their path around the circle, and may thus repeat a_1 rotations as many times as necessary until they land in I'_1 . In conclusion, the number of times we need to rotate to come closer to 1 than $a_1\alpha$ for the first time is of the form $n \cdot a_1\alpha + 1$, where $n \geq 1$ is an integer.

If we in I'_0 and I'_1 for illustrative purposes substitute 1 for 0, $a_1\alpha$ for α , and $1 - \alpha$ for 1 we should find ourselves in a very familiar situation. We get $I_1 \sim [1, 0)$, $I'_0 \sim [1 - \alpha, 0)$, and $I'_1 \sim [1, 1 - \alpha)$. Ignoring the reversed orientation of the intervals, this is the same setup as we started with if we take I_1 as the entire circle, I'_0 as I_0 and I'_1 as the new I_1 .

From this point of view, it makes sense to make a renormalisation of I_1 , mapping it from $[1 - \alpha, 1)$ back to S^1 , i.e. the circle $[0, 1)$. By doing so we can then recast the problem of finding the next closest point to 0 in the sequence to the problem of finding α_1 . Define a new rotation map for the renormalised circle. We want the new rotation map to correspond to the rotation $\alpha_1\alpha$ in the old circle, and so find $\alpha_2 = \frac{1-\alpha_1\alpha}{\alpha}$, where we have reversed the rotation so it is the same direction as in the original. Our new intervals become $I'_0 = [0, 1 - \alpha_2)$, and $I'_1 = [1 - \alpha_2, 1)$. Again, we write the sequence up to and including the first 1 in this new rotation as $C_2 = 0^{\alpha_2-1}1$. We can then translate this back to symbols in our original sequence by noting that each zero in C_2 corresponds to a C_1 in σ . Then we need to make one more α_2 -corresponding rotation in the original circle to move out of I_1 , corresponding to 0^{α_1} in σ , and a final α rotation giving 1. These last two steps can be rewritten as $0^{\alpha_1}1 = 00^{\alpha_1-1}1 = 0C_1$. To summarize, the first time the sequence gets closer to 1 than $\alpha_1\alpha$, is after C_1C_2 , where we have denoted the added steps after reaching $\alpha_1\alpha$ as $C_2 = C_1^{\alpha_2-1}0C_1$.

It should not be an unreasonable proposition that we can iteratively repeat the renormalisation inside each most recent renormalised circle, getting us closer and closer to 1, giving a sequence of C_n . In that case, $\sigma = C_1C_2C_3\dots$, and we will find generally that

$$\begin{aligned} C_1 &= 0^{\alpha_1-1}1, \\ C_2 &= C_1^{\alpha_2-1}0C_1, \\ &\vdots \\ C_n &= C_{n-1}^{\alpha_n-1}C_{n-2}C_{n-1} \end{aligned}$$

Recall that the α_n were found by making the smallest number of rotations of length α that were needed to get closer to 0. We have also related these rotations to diophantine equations, and we may through this connection view the rotation sequence as a visualisation of the euclidean algorithm. Specifically, the renormalisations correspond to shifting the focus to the rest, and the rotations to finding the integer part. Because of this, the α_n are exactly the continued fraction coefficients for α .

Theorem 4.0.6. *All rotation sequences are Sturmian and every Sturmian sequence can be uniquely encoded in rotation sequences by choice of an irrational number α .*

Proof. We need to prove that any rotation sequence, R , has the correct complexity function. Let us first simply remark that there are two words of length 1 in the rotation sequence, namely 0 or 1. As inductive step, assume there are $n + 1$ unique words of length n in R . Then we need to show that there are exactly $n + 2$ words of length $n + 1$.

Our strategy will depend on the fact that every word of length n also appears in a word of length $n + 1$, but with either an extra zero or one at the end. Thus we need to demonstrate that exactly one of the words of length n appears twice among the longer words, in one case extended by zero and in another by one. This word cannot end in a 1, since a point in I_1 always ends up in I_0 after one rotation. By similar reasoning, we may also conclude that the word must end in a neighbourhood of $1 - 2\alpha$, as that is the only set of point for which the next rotation straddles the boundary between I_0 and I_1 . We have previously shown that the rotation sequence gets arbitrarily close to 0 on both sides, and it therefore also gets arbitrarily close on either side of $1 - 2\alpha$, two rotations earlier.

The last step of the proof is to substantiate the claim that words ending in a sufficiently small neighbourhood of $1 - 2\alpha$ are the same, or equivalently, small changes of the endpoint (or starting point) on the circle of such words does not change the sequence. This is true, since otherwise a previous point in the rotation sequence crosses between I_0 and I_1 at the exact same time as the endpoint crosses $1 - 2\alpha$. Let us instead assume such a point exists, which we can write as $p = 0 = k\alpha + \beta$, for some choice of β and $0 < k < n$. This implies $(k + (n - k))\alpha + \beta = n\alpha + \beta = 1 - 2\alpha$ is reachable from p . Then $(n + 2)\alpha + \beta = 0 = p$. But this means that the rotation has gone from p back to itself, and therefore α must be rational, which is a contradiction. So in conclusion there is a neighbourhood containing $1 - 2\alpha$ such that words ending there are the same. But words ending before $1 - 2\alpha$ will next rotation append a 0 to the rotation sequence, and words ending after will append a 1. Since this neighbourhood is unique in this property, only one word of length n will be extended by both 0 and 1, and the sequence as a whole is therefore Sturmian. \square

Additional Notes

A variant of the rotation sequence is the cutting sequence in the euclidean plane.

Definition 4.0.7. A cutting sequence is a sequence built by taking a line, $y = \alpha x + \beta$, on an integer grid with the line going from the origin to positive infinity. For each vertical gridline the line intersects (cuts), a 0 is appended to the sequence, and for each horizontal line, a 1. If $\alpha + \beta$ is rational, the line might at some point cut a horizontal and a vertical gridline at the same time, at which point either 0 or 1 may be appended at will.

These kinds of sequences are also Sturmian, and it is therefore possible to transform the α and β between the cutting sequence and rotation sequence such that both yield the same Sturmian word. Specifically, let us distinguish the rotation sequence parameters by writing

α_R and β_R . The corresponding cutting sequence then has parameters $\alpha_C = \frac{\alpha_R}{1+\alpha_R}$ and $\beta_C = \frac{\beta_R}{1+\alpha_R}$. Since α_C and α_R is not the same, it is important to mind which representation is used when referring to a Sturmian sequence by only α .

A concrete example of a Sturmian sequence is the Fibonacci word, which is the word corresponding to the cutting sequence with $\alpha = 1/\phi^2$. The Fibonacci word is formed by letting $S_0 = 0$ and $S_1 = 01$, $S_n = S_{n-1}S_{n-2}$ and then letting n go to infinity.

As mentioned earlier, Sturmian sequences also occur in other contexts, and one such variant is the Beatty sequences.

Definition 4.0.8. A *Beatty sequence*, denoted \mathcal{B} , is a sequence formed by taking any positive irrational number, θ , and writing down the integer parts of its integer multiples, plus some initial value given by $\gamma \in [0, 1)$. Algebraically this can be expressed by

$$\mathcal{B}_\theta = (\lfloor \theta + \gamma \rfloor, \lfloor 2\theta + \gamma \rfloor, \lfloor 3\theta + \gamma \rfloor, \dots).$$

Sturmian sequences are the characteristic sequences of Beatty sequences, in the sense that the Sturmian terms indicate the difference between consecutive Beatty terms. Specifically,

$$\lfloor (n+1)\theta + \gamma \rfloor - \lfloor n\theta + \gamma \rfloor,$$

can only take one of two values, $\lfloor \theta \rfloor$ or $\lfloor \theta \rfloor + 1$. In other words the differences form a binary sequence, which we claim is Sturmian. This is easy to show for $\theta \in [0, 1]$, since the difference forms exactly a rotation sequence, and by noting that the ordering of the symbols only depends on the fractional part, the claim must be true for other θ as well. Beatty sequences have other interesting properties, for example if θ is combined with another irrational number ψ , such that $\frac{1}{\theta} + \frac{1}{\psi} = 1$, then the Beatty sequences for both numbers contain every integer exactly once.

HYPERBOLIC GEOMETRY

The next system we shall explore is found in the tessellation of the hyperbolic plane. To understand what this means we will need some basic knowledge of what hyperbolic geometry is, starting with how to model the hyperbolic plane. Locally we can think of it as similar to the ordinary euclidean plane. Globally however, there is more space inbetween any two points; one might imagine two persons moving in different directions from the same starting point. They would experience the distance between them increasing exponentially, as opposed to linearly. More formally, given a surface, M , and a metric $ds^2 = E(x, y)dx^2 + F(x, y)dxdy + G(x, y)dy^2$, M is said to be *hyperbolic* if it has constant negative *Gaussian curvature*, K , defined by

$$K = -\frac{1}{2\sqrt{EG}} \left(\frac{\delta}{\delta x} \frac{E'_x}{\sqrt{EG}} + \frac{\delta}{\delta y} \frac{G'_y}{\sqrt{EG}} \right).$$

All two-dimensional metrics may be expressed by the above form, with constant curvatures typically normalised to 1, 0, or -1 , depending on sign. Another characterisation of E and G can be found by letting $X(x, y)$ be some orthogonal parametrisation of M . Then we can calculate $E = (X'_x, X'_x)$ and $G = (X'_y, X'_y)$. In other words, x and y can be thought of as a coordinate system over M , from which the tangent plane spanning vectors X'_x and X'_y are then involved in calculation of the curvature. A more indepth explanation on the topic can be found in the lecture notes by Sigmundur Gudmundsson [9].

Since the hyperbolic plane is difficult to embed in euclidean space, we will begin with a model of it in the form of the upper half plane, $\mathbb{H} = \{z \in \mathbb{C} : \text{Im}(z) > 0\}$, also called the Lobachevsky plane. To measure distance in \mathbb{H} we shall use the metric $ds = \frac{|dz|}{\text{Im}(z)} = \frac{\sqrt{dx^2 + dy^2}}{y}$. The length of a path γ is then defined by $l(\gamma) = \int_{\gamma} ds$. By parametrising as $\gamma(t) = x(t) + iy(t)$, $t \in [0, 1]$, we have

$$l(\gamma) = \int_0^1 \frac{1}{y(t)} \sqrt{\left(\frac{dx}{dt}\right)^2 + \left(\frac{dy}{dt}\right)^2} dt.$$

We define the distance between two points $p, q \in \mathbb{H}$ as

$$d_{\mathbb{H}}(p, q) = \inf \{l(\gamma) : \gamma(0) = p, \gamma(1) = q\}.$$

We claim that $d_{\mathbb{H}}(p, q)$ is a metric.

Proof. Firstly, $l(\gamma)$ is an integral over positive values, so $d_{\mathbb{H}}(p, q) \geq 0$. Because $\text{Im}(z) = y > 0$, if $d_{\mathbb{H}}(p, q) = 0$ the parametrised form shows

that there must be no change of x or y when t varies, i.e. $\frac{dx}{dt} = \frac{dy}{dt} = 0$. Hence $\gamma(0) = \gamma(1)$, so a distance of zero implies $p = q$. Since the integral cannot be smaller than 0, $p = q$ implies a distance of zero.

Secondly, it is possible to substitute t for $u = 1 - t$ in the parametrised integral, which gives

$$\begin{aligned} l(\gamma(t)) &= \int_0^1 \frac{1}{y(t)} \sqrt{\left(\frac{dx}{dt}\right)^2 + \left(\frac{dy}{dt}\right)^2} dt \\ &= \int_1^0 \frac{1}{y(u)} \sqrt{\left(\frac{dx}{du}\right)^2 + \left(\frac{dy}{du}\right)^2} du \\ &= \int_0^1 \frac{1}{y(1-t)} \sqrt{\left(\frac{dx}{dt}\right)^2 + \left(\frac{dy}{dt}\right)^2} dt = l(\gamma(1-t)). \end{aligned}$$

The set of all paths from p to q are the same paths, but with reversed direction, going from q to p , meaning $d_{\mathbb{H}}(p, q) = d_{\mathbb{H}}(q, p)$.

Thirdly, the triangle equality holds. Suppose in contradictory terms that $d_{\mathbb{H}}(p, r) > d_{\mathbb{H}}(p, q) + d_{\mathbb{H}}(q, r)$. Then there exists at least two paths, γ_1 and γ_2 , such that $d_{\mathbb{H}}(p, q) = l(\gamma_1)$ and $d_{\mathbb{H}}(q, r) = l(\gamma_2)$. Thus the concatenation of γ_1 and γ_2 , call it $\gamma_1 + \gamma_2$, connects p and r , but we also have $l(\gamma_1 + \gamma_2) = d_{\mathbb{H}}(p, q) + d_{\mathbb{H}}(q, r)$. This contradicts the initial supposition, and so $d_{\mathbb{H}}(p, r) \leq d_{\mathbb{H}}(p, q) + d_{\mathbb{H}}(q, r)$.

With that, all conditions of a metric have been shown to hold. \square

Now that we have a metric at hand we are prepared to check that \mathbb{H} is an accurate representation the hyperbolic plane, by confirming that the curvature is -1 as expected. First we identify $E = G = \frac{1}{y^2}$. We will also need the partial derivatives, $E'_x = 0$, $G'_y = -\frac{2}{y^3}$. The gaussian curvature is then

$$\kappa = -\frac{\sqrt{y^4}}{2} \left(\frac{\delta}{\delta x} (0 \cdot \sqrt{y^4}) - \frac{\delta}{\delta y} \frac{2\sqrt{y^4}}{y^3} \right) = -\frac{y^2}{2} \cdot \frac{2}{y^2} = -1.$$

So \mathbb{H} has a constant negative curvature! With that we will introduce a characterisation of lines that is workable in the context of curved surfaces.

Definition 5.0.1. A path is called *geodesic* if it locally minimises distances. This can be interpreted as a path traced by only moving forwards.

In hyperbolic space, minimising distance locally is the same as minimising distance globally, so a geodesic is a shortest path between points. Contrast this with the geodesics on a torus, where there are infinitely many geodesics connecting two points, but only one minimizing distance.

Just as in euclidean geometry, linear transformations serve an important role in both defining and working in the hyperbolic plane. As such, the following definition introduces one of (if not the most) important concepts in hyperbolic maps.

Definition 5.0.2. A *Möbius map* is a map of the form

$$f : \hat{\mathbb{C}} \rightarrow \hat{\mathbb{C}}, z \mapsto \frac{az + b}{cz + d},$$

where $a, b, c, d \in \mathbb{R}$ and $ad - bc \neq 0$.

Here we also define $\hat{\mathbb{C}} = \mathbb{C} \cup \infty$, where $f(\infty) = \lim_{|z| \rightarrow \infty} f(z)$, and $\frac{z}{0} = \infty$ when $z \notin \{0, \infty\}$. For z zero or infinity, we leave $\frac{z}{0}$ undefined. We note in particular that

$$f(\infty) = \lim_{|z| \rightarrow \infty} \frac{az + b}{cz + d} = \frac{a}{c},$$

which always holds, since in the case a and c are 0, then $ad - bc = 0$, and as such we are not working with a Möbius map. We also have

$$f\left(-\frac{d}{c}\right) = \frac{-a \cdot \frac{d}{c} + b}{0} = \infty,$$

which is well defined because $-a \cdot \frac{d}{c} + b = 0$ also implies that $ad - bc = 0$.

Möbius maps are particularly convenient to work with as they fulfill certain helpful properties.

Definition 5.0.3. The conformal (angle preserving) bijections of $\hat{\mathbb{C}}$ are denoted $\text{Aut}(\hat{\mathbb{C}})$, the automorphism group of $\hat{\mathbb{C}}$. Automorphisms are angle preserving, fully invertible, structure-preserving maps (isomorphisms) from a space to itself.

Theorem 5.0.4. *Möbius maps are bijective and conformal. In other words, Möbius maps make out a subset of $\text{Aut}(\hat{\mathbb{C}})$.*

Proof. Möbius maps are invertible in $\hat{\mathbb{C}}$ by $f^{-1}(z) = \frac{dz - b}{-cz + a}$, since

$$f(f^{-1}(z)) = \frac{a \frac{dz - b}{-cz + a} + b}{c \frac{dz - b}{-cz + a} + d} = \frac{adz - ab - bcz + ab}{cdz - bc - cdz + ad} = \frac{(ad - bc)z}{ad - bc} = z.$$

Note that $f^{-1}\left(\frac{a}{c}\right)$ is well defined because $d = \frac{bc}{a} \Rightarrow ad - bc = 0$. We also confirm that

$$f^{-1}\left(f\left(-\frac{d}{c}\right)\right) = f^{-1}(\infty) = -\frac{d}{c},$$

and

$$f^{-1}(f(\infty)) = f^{-1}\left(\frac{a}{c}\right) = \frac{\frac{da}{c} - b}{0} = \infty.$$

This implies that the maps are bijective.

The derivative of a Möbius map is

$$f' = \frac{ad - bc}{(cz + d)^2},$$

which, since $ad - bc \neq 0$, shows that the map is complex differentiable everywhere except at 0. That it is conformal then follows from it being holomorphic. \square

From the proof we also get an important second result that the inverse of a Möbius map is also a Möbius map.

It might have occurred to observant readers that the condition $ad - bc \neq 0$ is strangely reminiscent of determinants from linear algebra. Indeed, there is more than a superficial connection between the two.

Theorem 5.0.5. *The map*

$$F : GL_2(\mathbb{C}) \rightarrow \text{Aut}(\hat{\mathbb{C}}), \begin{pmatrix} a & b \\ c & d \end{pmatrix} \mapsto \left(f(z) = \frac{az + b}{cz + d} \right)$$

is a homomorphism with respect to matrix multiplication and function composition, where the set GL_2 is the (general) group of linearly independent 2×2 -matrices. This means that the composition of two Möbius maps can be interpreted as a matrix multiplication.

Proof. Assume we have two Möbius maps, $f(z) = \frac{az + b}{cz + d}$, and $g(z) = \frac{sz + t}{uz + v}$. Then let

$$A = \begin{pmatrix} s & t \\ u & v \end{pmatrix} \begin{pmatrix} a & b \\ c & d \end{pmatrix} = \begin{pmatrix} as + ct & bs + dt \\ au + cv & bu + dv \end{pmatrix}.$$

We also have

$$\begin{aligned} (f \circ g)(z) &= \frac{s \frac{az+b}{cz+d} + t}{u \frac{az+b}{cz+d} + v} \\ &= \frac{asz + sb + ctz + td}{auz + ub + cvz + dv} = \frac{(as + ct)z + bs + dt}{(au + cv)z + bu + dv}, \end{aligned}$$

which makes it clear that $F(A) = (f \circ g)$ □

As defined, F takes matrices from $GL_2(\mathbb{C})$ to Möbius maps, and had Möbius maps been a strict subset of $\text{Aut}(\hat{\mathbb{C}})$ it might have been reasonable to introduce notation to denote the set of all Möbius transformations, so as to make F a surjection. However, the automorphisms of $\hat{\mathbb{C}}$ and the Möbius transformations are one and the same.

Theorem 5.0.6. *The map F is surjective.*

We will leave this proof unsolved in order to not stray too far from our intended discourse on hyperbolic geometry, as it requires theory from analysis.

The crux of this theorem is in the given type of F , namely it is from matrices to $\text{Aut}(\hat{\mathbb{C}})$ rather than to Möbius transformations. To put it in other words, all transformations in $\text{Aut}(\hat{\mathbb{C}})$ can be naturally expressed as matrices, but there is some ambiguity preventing F from being a bijection, since $\begin{pmatrix} ka & kb \\ kc & kd \end{pmatrix}$ and $\begin{pmatrix} a & b \\ c & d \end{pmatrix}$ map to the same element in

$\text{Aut}(\hat{\mathbb{C}})$. Instead we can restrict our use of matrices to those with a determinant of ± 1 , i.e.

$$\text{SL}_2(\mathbb{C}) = \{A \in \text{GL}_2(\mathbb{C}) : \det(A) = 1\}.$$

This still allows both A and $-A$, so we form the quotient set

$$\text{PSL}_2(\mathbb{C}) = \text{SL}_2(\mathbb{C}) / \pm I,$$

meaning multiplication by $-I$ gives members of the same equivalence class. In simpler terms, A and $-A$ are considered to be the same matrix in $\text{PSL}_2(\mathbb{C})$. By only allowing matrices from $\text{PSL}_2(\mathbb{C})$, F becomes an isomorphism from such matrices to $\text{Aut}(\hat{\mathbb{C}})$, and we will from here on out use matrices to denote members of $\text{Aut}(\hat{\mathbb{C}})$ interchangeably with function notation.

Theorem 5.0.7. *Möbius transformations can be decomposed into three basic maps:*

TRANSLATION $f(z) = \begin{pmatrix} 1 & \alpha \\ 0 & 1 \end{pmatrix} = z + \alpha$ where $\alpha \in \mathbb{R}$,

DILATION $f(z) = \begin{pmatrix} \beta & 0 \\ 0 & \frac{1}{\beta} \end{pmatrix} = \beta^2 z$, where $\beta \in \mathbb{R}^+$,

INVERSION $f(z) = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} = -1/z$.

Proof. Suppose $T = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$. Then define the following

$$T_1 = \begin{pmatrix} c & 0 \\ 0 & \frac{1}{c} \end{pmatrix}, T_2 = \begin{pmatrix} 1 & \frac{d}{c} \\ 0 & 1 \end{pmatrix},$$

$$T_3 = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}, T_4 = \begin{pmatrix} 1 & \frac{a}{c} \\ 0 & 1 \end{pmatrix}.$$

Note in particular that all of them are variants of the three basic maps defined in the theorem statement. Then,

$$\begin{aligned} T_4 T_3 T_2 T_1 &= T_4 T_3 \begin{pmatrix} c & d \\ 0 & \frac{1}{c} \end{pmatrix} \\ &= T_4 \begin{pmatrix} 0 & -\frac{1}{c} \\ c & d \end{pmatrix} \\ &= \begin{pmatrix} \frac{ca}{c} & \frac{ad}{c} - \frac{1}{c} \\ c & d \end{pmatrix} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}. \end{aligned}$$

In the last step we made use of $ad - bc = 1$, implying $b = \frac{ad-1}{c}$. \square

For the sake of completeness, it should be mentioned that Möbius maps are typically defined in a more general manner, where α and β are allowed to take complex values. In particular, $\text{Aut}(\hat{\mathbb{C}})$ is often defined to allow the coefficients a, b, c, d to be complex. The proof given here should still hold for that case. However, as we now move on to demonstrating the relationship between geodesics and Möbius transformations in \mathbb{H} , we must require the coefficients to be real.

Theorem 5.0.8. *The set consisting of vertical lines contained in \mathbb{H} and semicircles in centered on \mathbb{R} contained in \mathbb{H} is closed under $T \in \text{Aut}(\hat{\mathbb{C}})$.*

Proof. We only need to show this holds for each of the three basic maps from the previous theorem. The statement is trivially true for translation and dilation.

Then remains inversion. As a side note, in the theorem statement we are not entirely specific on the conditions for when a line should go to a circle and when it should go to a line, and so it is possible to prove the theorem in this case in a relatively simple manner. We shall nonetheless take a longer route since it is more instructive.

A vertical line in \mathbb{H} may be described by $z = c + iy$, where c is constant, and $y \in [0, \infty]$. Then

$$-\frac{1}{z} = -\frac{1}{c + iy} = \frac{-c + iy}{c^2 + y^2}.$$

Now we perform the variable substitution

$$y = \sqrt{\left(\frac{1}{\sin^2 t} - 1\right) c^2},$$

such that $t \in [0, \pi/2]$, giving

$$\frac{1}{c^2 + y^2} = \frac{1}{c^2 + \frac{c^2}{\sin^2 t} - c^2} = \frac{\sin^2 t}{c^2}.$$

Furthermore,

$$\begin{aligned} -\frac{1}{z} &= \frac{-c + iy}{c^2 + y^2} \\ &= -\frac{\sin^2 t}{c} + i \frac{\sqrt{\frac{1}{\sin^2 t} - 1} |c|}{c^2} \sin^2 t \\ &= \frac{\cos(2t) - 1}{2c} + i \frac{\sqrt{(1 - \sin^2 t) \sin^2 t}}{|c|} \\ &= \frac{\cos(2t) - 1}{2c} + i \frac{\cos t \sin t}{|c|} = \frac{\cos(2t) - 1}{2c} + i \frac{\sin(2t)}{2|c|} \end{aligned}$$

Assuming $c \neq 0$, this describes a semicircle contained in \mathbb{H} , as the imaginary part is never negative, and $t = 0$ gives $-\frac{1}{z} = 0$ (the origin), and $t = \pi/2$ gives $-\frac{1}{z} = -\frac{1}{2c}$, meaning the endpoints lie on $\mathbb{R} \cup \{\infty\}$. When $c = 0$, we have the imaginary axis, for which $-\frac{1}{z} = \frac{i}{y}$, that is, the imaginary axis maps to itself.

It remains to show that circles are mapped to circles. The equation for a circle centered at α in the complex plane can be described by

$$(z - \alpha)(\bar{z} - \bar{\alpha}) = z\bar{z} - \alpha\bar{z} - \bar{\alpha}z + \alpha\bar{\alpha} = r^2.$$

Then if $T(z) = \omega = -\frac{1}{z}$, we get $z = -\frac{1}{\omega}$, and

$$\frac{1}{\omega\bar{\omega}} - \frac{\alpha}{\bar{\omega}} - \frac{\bar{\alpha}}{\omega} + \alpha\bar{\alpha} = r^2.$$

Assuming $\alpha \neq r$,

$$1 - \alpha\omega - \bar{\alpha}\bar{\omega} + (|\alpha|^2 - r^2)|\omega|^2 = 0 \Rightarrow \frac{|\alpha|^2 - r^2}{(|\alpha|^2 - r^2)^2} - \frac{\alpha\omega + \bar{\alpha}\bar{\omega}}{|\alpha|^2 - r^2} + |\omega|^2 = 0$$

Using $\alpha\omega + \bar{\alpha}\bar{\omega} = 2\operatorname{Re}(\omega\bar{\alpha})$, we can rewrite this is in the form of a completed the square, giving

$$\left| \omega - \frac{\bar{\alpha}}{|\alpha|^2 - r^2} \right|^2 = \frac{r^2}{|\alpha|^2 - r^2},$$

which is the equation for a circle. When instead $|\alpha| = r$, the circle cuts through the origin. As noted in the case for the line, any vertical line maps to such a circle, and since $T(T(z)) = z$, circles cutting the origin maps to a vertical line. \square

Again for the sake of completeness, the same theorem can be stated in a more general manner when working with all of \mathbb{C} (rather than only \mathbb{H}) and complex Möbius map coefficients. Then the theorem states that any line and circle maps to either a line or circle.

Let us now turn back to \mathbb{H} , where we will take one more step towards finding the geodesics.

Theorem 5.0.9. *The maps in $\operatorname{Aut}(\hat{\mathbb{C}})$ act as isometries on \mathbb{H} .*

Proof. Let $T = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$, and $w = T(z) = \frac{az+b}{cz+d}$. Then

$$\frac{dw}{dz} = \frac{a(cz+d) - (az+b)c}{(cz+d)^2} = \frac{ad-bc}{(cz+d)^2} = \frac{1}{(cz+d)^2},$$

meaning $dw = \frac{dz}{(cz+d)^2}$. Also, writing $z = x + iy$ gives

$$\begin{aligned} w &= \frac{ax + iay + b}{cx + icy + d} \\ &= \frac{(ax + b + iay)(cx + d - icy)}{(cx + d + icy)(cx + d - icy)} \\ &= \frac{ac(x^2 + y^2) + x + iy + bd}{(cx + d + icy)(cx + d - icy)} \\ &= \frac{ac|z|^2 + z + bd}{|cz + d|^2} \end{aligned}$$

Therefore $\text{Im}(w) = \frac{\text{Im}(z)}{|cz+d|^2}$. By the definition of curve length in \mathbb{H} ,

$$l(T(\gamma)) = \int_{T(\gamma)} \frac{|dw|}{\text{Im}(w)} = \int_{\gamma} \frac{|cz + d|^2}{\text{Im}(z)} \cdot \frac{|dz|}{|cz + d|^2} = \int_{\gamma} \frac{|dz|}{\text{Im}(z)} = l(\gamma).$$

□

The last theorem before we are ready to deal with the geodesics explains the automorphisms, $\text{Aut}(\mathbb{H}) = \{T \in \text{Aut}(\hat{\mathbb{C}}) : T(\mathbb{H}) = \mathbb{H}\}$.

Theorem 5.0.10. *The automorphisms of \mathbb{H} are $\text{Aut}(\hat{\mathbb{C}})$.*

Proof. Up until this point we have only used \mathbb{H} as a model for the hyperbolic plane. Another very common model is the Poincaré disk,

$$\mathbb{D} = \{z \in \mathbb{C} : |z| < 1\}.$$

It is easy to swap between \mathbb{H} and \mathbb{D} by using the Cayley transform,

$$f(z) = \frac{z - 1}{z + 1}.$$

We can see that the Cayley transform moves \mathbb{H} to \mathbb{D} since for example $f(0) = -1$, $f(1) = -i$ and $f(\infty) = 1$. Now, since $\mathbb{R} \cup \infty$ can be considered a circle with infinite radius, and since Möbius transformations map circles to circles, the Cayley transform must map \mathbb{R} to the unit circle. In a similar fashion $f(i) = 0$, so the rest of \mathbb{H} goes inside the unit disk, as opposed to around it. With the help of the Cayley transform we can then apply all we know in \mathbb{H} on \mathbb{D} and vice versa.

According to Schwarz Lemma, a holomorphic map $f : \mathbb{D} \rightarrow \mathbb{D}$ is an automorphism if it is isometric with respect to the Poincaré metric, which in turn can be thought of as the metric in \mathbb{D} keeping distances in \mathbb{H} constant under the Cayley transform. We know that maps in $\text{Aut}(\mathbb{H})$ are isometric, so if we take maps from $\text{Aut}(\hat{\mathbb{C}})$ we can compose them with the Cayley transform to find isometric, and therefore automorphic maps in \mathbb{D} . Möbius maps are fully invertible, so we may use the inverse Cayley transform to show that we also have the automorphism property in \mathbb{H} . □

Theorem 5.0.11. *Geodesics in \mathbb{H} are either vertical lines or circular arcs centered on the real axis.*

Proof. First let us consider the case of vertical lines. Consider two points on the imaginary axis, ai and bi such that $b > a$. The vertical path connecting them, $\gamma_{ab} = iy$ for $y \in [a, b]$, has length

$$l(\gamma_{ab}) = \int_a^b \frac{1}{y} \sqrt{dx^2 + dy^2} = \int_a^b \frac{1}{y} dy = \log \frac{b}{a}.$$

Any deviation from the vertical path will simply add to the integral, as $\frac{1}{y}$ is positive and constant with respect to x .

Now let us consider circular arcs. We know that the positive imaginary axis is a geodesic, and any Möbius mapping of the imaginary axis will also be a semicircle (or line). Möbius maps are isometries, so any map of \mathcal{I} will also be a geodesic. Semicircles and lines are closed under these transformations, and by inversion, translation and scaling we can reach any circular arc centered on \mathbb{R} using \mathcal{I} .

Finally, all isometries in \mathbb{H} are Möbius maps. Therefore, all geodesics must be expressible as a transformation of the imaginary axis through a Möbius transformation, which we have just described. Hence there are no geodesics other than those vertical lines and semicircles. \square

As geodesics are essentially a generalisation of lines, we are now able to discuss tilings of \mathbb{H} . Analogous to the grid in the euclidean plane for the cutting sequence, let us introduce such a tiling, called the Farey tessellation.

Definition 5.0.12. The *Farey tessellation*, denoted \mathcal{F} , is a tiling of the hyperbolic plane, and can be constructed in \mathbb{H} with the following steps.

1. Start with the set \mathcal{F}_0 of geodesics connecting $n \in \mathbb{Z}$ to ∞ . These are vertical lines. Label each integer n by $\frac{n}{1}$, and call adjacent integers *neighbours*.
2. Connect each pair of neighbours with the geodesic that has both neighbours as an endpoint. Add these geodesics to \mathcal{F}_{k-1} , creating \mathcal{F}_k .
3. For each pair of neighbours $(\frac{p}{q}, \frac{r}{s})$, mutually replace the other neighbour with the number at $\frac{p}{q} \oplus \frac{r}{s} := \frac{p+r}{q+s}$ as a new neighbour. Take for example $\frac{1}{1}$ and $\frac{2}{1}$. Then $\frac{1}{1}$ replaces $\frac{2}{1}$ by $\frac{1+2}{1+1} = \frac{3}{2}$ as its new neighbour, and $\frac{2}{1}$ replaces $\frac{1}{1}$ by $\frac{3}{2}$. The fraction $\frac{p+r}{q+s}$ should be considered both a position on \mathbb{R} and a label.
4. Repeat steps 2-3.

The tiling is then $\mathcal{F} = \lim_{k \rightarrow \infty} \mathcal{F}_k$, as seen in Figure 3, where it is illustrated after four iterations.

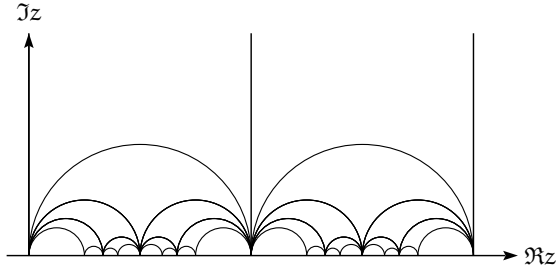


Figure 3: The Farey tessellation after four iterations, for values $\text{Re}(z) \in [0, 2]$.

Before we prove that \mathcal{F} actually tiles \mathbb{H} , let us discuss how we should look at it intuitively. Let us make clear what the tiles are. Each time at step 3 in the construction, we take a pair of neighbours, $(\frac{p}{q}, \frac{r}{s})$ and make pairs $(\frac{p}{q}, \frac{p+r}{q+s})$ and $(\frac{p+r}{q+s}, \frac{r}{s})$. After we return to step 2, all three pairs are connected by a geodesic, and it is these that together forms hyperbolic equivalent of a triangle since the geodesics pairwise share one endpoint (label), and no endpoint is shared between all three. Note also that the first time we reach step 2, we can form triangles from \mathcal{F}_1 , namely $(\frac{n}{1}, \frac{n+1}{1}, \infty)$. As such, the tiles of \mathcal{F} are *triangles*. Since all corners lie on $\mathbb{R} \cup \infty$, which corresponds to a circle with infinite radius in the euclidean plane, we call the triangles *ideal*.

For the triangles to be tiles of the Farey tessellation they collectively need to fill up \mathbb{H} entirely, without overlaps. Because of this, each triangle can be taken as a fundamental domain under $\text{PSL}_2(\mathbb{Z})$, as we will show in the following theorem where we will use the triangle $\Delta = (0, 1, \infty)$.

Theorem 5.0.13. *The farey tiling \mathcal{F} tiles \mathbb{H} , and the automorphisms of \mathcal{F} are $\text{Aut}(\mathcal{F}) = \text{PSL}_2(\mathbb{Z})$.*

Proof. Consider $T = \begin{pmatrix} r & p \\ s & q \end{pmatrix} \in \text{PSL}_2(\mathbb{Z})$ such that $\frac{p}{q} < \frac{r}{s}$ are connected by a geodesic in \mathcal{F} . Then we calculate $T(0) = \frac{p}{q}$, $T(\infty) = \frac{r}{s}$, and $T(1) = \frac{p+r}{q+s}$, which means the triangle $\Delta = (0, 1, \infty)$ maps to that of $T(\Delta) = (\frac{p}{q}, \frac{p+r}{q+s}, \frac{r}{s})$. Since $(\frac{p}{q}, \frac{r}{s})$ were chosen arbitrarily from \mathcal{F} , each triangle can be reached from any other by transforming to Δ with coefficients p_1, q_1, r_1, s_1 using T_1^{-1} , then from Δ using T_2 with coefficients p_2, q_2, r_2, s_2 . Members of $\text{PSL}_2(\mathbb{Z})$ are conformal bijections, and so we get $\text{Aut}(\mathcal{F}) \subset \text{PSL}_2(\mathbb{Z})$.

It should be clear from construction that \mathcal{F} covers \mathbb{H} , now we show that there are no overlaps. Let us assume the contrary, that two triangles, $t_1, t_2 \in \mathcal{F}$, overlap at some point. Then we are able to find a map $M_1 \in \text{PSL}_2(\mathbb{Z})$ taking for example t_1 to Δ , such that $M_1(t_2)$

cuts at least one of its sides. Let us introduce $S = \begin{pmatrix} 0 & -1 \\ 1 & -1 \end{pmatrix}$, for which $S(0) = 1$, $S(1) = \infty$, and $S(\infty) = 0$. In other words, S takes Δ to itself, but rotated, and we can use S to rotate $M_1(t_2)$ such that the side cut

by $M_1(t_2)$ is the imaginary axis. Since $\text{Aut}(\mathcal{F})$ is a subset of $\text{PSL}_2(\mathbb{Z})$, there must also be another map $M_2 = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \in \text{PSL}_2(\mathbb{Z})$, such that $M_2(\Delta) = M_1(t_2)$. In particular, one of the sides that cuts the imaginary axis has endpoints in $\frac{a}{c}$ and $\frac{b}{d}$. Consider the case $\frac{b}{d} < 0 < \frac{a}{c}$, where we also take $d > 0$, factoring out a -1 if necessary. Then b is negative, and a and c have the same sign. Additionally a, b, c, d are non-zero integers, thus we have $|ad - bc| \geq 2$, and M_2 could not be in $\text{PSL}_2(\mathbb{Z})$. If $\frac{a}{c} < 0 < \frac{b}{d}$, we may take $c > 0$ to achieve the same result. Therefore our assumption that t_1 and t_2 overlap must be false.

Now no member of $\text{PSL}_2(\mathbb{Z})$ can map Δ to only partially overlap with another triangle of \mathcal{F} . It follows that $\text{PSL}_2(\mathbb{Z}) \subset \text{Aut}(\mathcal{F})$.

Above we have assumed that if $\frac{p}{q}$ and $\frac{r}{s}$ are neighbours, then $T = \begin{pmatrix} r & p \\ s & q \end{pmatrix}$ is a element in $\text{PSL}_2(\mathbb{Z})$, which still remains to be

proven. Note that since we constructed \mathcal{F} from neighbours $(\frac{n}{1}, \frac{n+1}{1})$, the determinant of the corresponding map T is $n+1 - n = 1$. Now assume inductively that $rq - ps = 1$. We have for $(\frac{p}{q}, \frac{p+r}{q+s})$ the determinant $(p+r)q - p(q+s) = pq + rq - pq - ps = rq - ps = 1$. Likewise for the neighbour pair $(\frac{p+r}{q+s}, \frac{r}{s})$, the determinant is $r(q+s) - s(p+r) = rq + rs - ps - rs = rq - ps = 1$. Hence any pair of points on \mathbb{R} connected by a geodesic in \mathcal{F} gives T a determinant of 1. \square

A useful property of \mathcal{F} follows from the above.

Theorem 5.0.14. *Every rational number $\frac{p}{q}$ corresponds to one endpoint of a geodesic in \mathcal{F} .*

Proof. By Bézout's theorem, there exists integers r, s such that $ps - qr = 1$, meaning we can define $T = \begin{pmatrix} p & r \\ q & s \end{pmatrix} \in \text{Aut}(\mathcal{F})$. For example \mathcal{J} has one endpoint in ∞ , thus $T(\infty) = \frac{p}{q}$ is also an endpoint. \square

Similar to the case for $\text{Aut}(\hat{\mathbb{C}})$, the automorphisms of \mathcal{F} can be fully described by only a small number of matrices.

Theorem 5.0.15. *Maps $P = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$ and $J = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}$ generate $\text{Aut}(\mathcal{F})$.*

Proof. Take an arbitrary matrix $T = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \in \text{Aut}(\mathcal{F})$, and consider how it composes with P and J . In the case of P ,

$$\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} a & b \\ c & d \end{pmatrix} = \begin{pmatrix} a+c & b+d \\ c & d \end{pmatrix}.$$

It is not difficult to see that repeated multiplication of P from the left gives

$$P^n T = \begin{pmatrix} a + nc & b + nd \\ c & d \end{pmatrix}.$$

The case of J is more simple, as J^2 is the identity matrix in $\text{PSL}_2(\mathbb{Z})$, and so only a single composition is of interest,

$$ST = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} a & b \\ c & d \end{pmatrix} = \begin{pmatrix} -c & -d \\ a & b \end{pmatrix}.$$

Now we are able to make steps similar to the euclidean algorithm, but with respect to the left-side elements of T . We write $\frac{a}{c} = n + \frac{r}{c}$, where n is the integer part, and $\frac{r}{c} \in [0, 1)$ the fractional part of the left hand side division. If $n \neq 0$, we write $a = nc + r$, and find that

$$P^{-n} T = \begin{pmatrix} a - nc & b - nd \\ c & d \end{pmatrix} = \begin{pmatrix} r & b - nd \\ c & d \end{pmatrix}.$$

The corresponding step in the euclidean algorithm would now be to perform the same steps for the division $\frac{c}{r}$, but to swap places between

r and c , we multiply from the left by J , thus $JP^{-n} T = \begin{pmatrix} -c & -d \\ r & b - nd \end{pmatrix}$.

Although the sign of c changed, the properties of the euclidean algorithm should still be applicable, meaning the rest r appearing in the upper left element of the matrix should eventually go to zero. A final swap of upper and lower elements by use of J gives something

of the form $T' = \begin{pmatrix} i & j \\ 0 & k \end{pmatrix}$. Now note that P^{-1} and J is in $\text{PSL}_2(\mathbb{Z})$,

meaning T' is too, and thus has determinant $ik = 1$. Moreover, i and k are integers, and therefore we have $i = k = \pm 1$, and we can factor

a -1 from T' to make them positive. Then $T' = \begin{pmatrix} 1 & j \\ 0 & 1 \end{pmatrix} = P^j$. Let us

call the composition of the matrices generated by P and J used in the euclidean-like algorithm S , and in particular note that $ST = P^n$ for some integer n . We have both P^n and S generated by P and J , and so $T = S^{-1}P^n$ is also generated by the same matrices. \square

This concludes the theory necessary theory for us to look back one last time to diophantine approximation and continued fractions.

Theorem 5.0.16. *Construct a cutting sequence for a real number $\alpha > 0$, by the following:*

- Start at an arbitrary point y_i on the positive imaginary axis \Im , in \mathbb{H} .

- Follow the oriented geodesic arc, γ , going from y_i to α . Any geodesic arc fulfilling this condition is sufficient.
- As the arc cuts through tiles of \mathcal{F} , note which sides are cut. Since the tiles are ideal triangles, exactly two sides of each triangle is cut. Should the sides meet on the right side of the oriented arc, append an R to the cutting sequence, otherwise an L.
- There is one exception to the pattern, which is when the geodesic intersects at a corner of a tile. This happens only on the real axis at α , at which point either L or R may be chosen, and the cutting sequence ends.

Following the above instructions will give a sequence of the form

$$L^{a_0} R^{a_1} L^{a_2} R^{a_3} \dots,$$

where for example L^3 means LLL, such that $a_k > 0$ for $k > 0$, and

$$[a_0; a_1, a_2, a_3, \dots] = \alpha.$$

Proof. Firstly, γ must cut a number of vertical lines equal to the integer part of α , after which it cuts its first semicircular arc in the interval $[a_0, a_0 + 1)$. The integer part is a_0 , so the sequence starts with $L^{a_0}R$.

Now we map the geodesic we cut to produce R back to \mathcal{J} , along with all of \mathcal{F} , thus resetting the oriented arc back to its starting point. The idea is that we can then read off the next time a non-vertical line will be cut in the same way as before. Generating matrices P and J preserve the orientation, meaning points left (or right) of γ are also to the left (right) of $P(\gamma)$ and $J(\gamma)$. In particular, P and J preserve the cutting sequence.

Let us mark the point at which γ cuts the last vertical line by z_0 . Then translation using P^{-a_0} takes z_0 to the imaginary axis, and $P^{-a_0}(\gamma)$ continues from $P^{-a_0}(z_0)$ to some point $P^{-a_0}(\alpha) = \alpha - a_0$, which lies between 0 and 1. Applying J, gives $JP^{-a_0}(z_0)$ still on the imaginary axis, and $JP^{-a_0}(\alpha) = -\frac{1}{\alpha - a_0}$. We know from the continued fraction representation of $\frac{1}{\alpha - a_0}$ that its integer part is a_1 . Thus JP^{-a_0} takes α to the interval $(-a_1 - 1, -a_1]$.

Since the sequence is preserved through the transformations, it remains to note that the imaginary axis represents the end of the sequence of Ls, and so our transformed γ will cut a_1 vertical lines, each giving R by the same argument as before. Then marking the last vertical line as z_1 , we can take it back to \mathcal{J} by P^{a_1} . Note specifically that we need to use a positive power this time, since the new α is negative. Then we can apply J again to find a_2 and so on. It should be clear that this will yield the cutting sequence as stated. \square

Additional notes

Related to the farey tessellation is the Stern-Brocot tree, which is an infinite binary tree with the property that each positive rational number corresponds to one of the vertices of the tree. Additionally, the root is labeled $\frac{1}{1}$, and each vertex $\frac{m}{n}$ with children $(\frac{p}{q}, \frac{r}{s})$ has that $\frac{p}{q} < \frac{r}{s}$, and $\frac{m}{n} = \frac{p}{q} \oplus \frac{r}{s} = \frac{p+r}{q+s}$. This is similar to the farey tessellation in that each (positive) rational number is represented in the structure, and the appearance of the operator \oplus . Particularly interesting is the fact that Stern-Brocot trees may be defined in terms of continued fractions, that is, for any fraction $\frac{p}{q} = [a_0; a_1, \dots, a_n - 1]$, its children are $[a_0; a_1, \dots, a_n + 1]$ and $[a_0; a_1, \dots, a_n - 1, 2]$ (recall that $[a_0; a_1, \dots, a_n, 1] = [a_0; a_1, \dots, a_n + 1]$, meaning $a_n \neq 1$). This suggests once again that continued fractions are deeply related to the farey tessellation.

Musing on the different systems where the continued fraction coefficients turn up, we have covered cases in two dimension with negative curvature, with the cutting sequence in the farey tessellation, as well as zero curvature, in the cutting sequence in the euclidean plane. Although maybe a bit of a stretch, the rotation sequence can be said to follow a geodesic on the surface of a sphere, which is a great circle. The unit sphere is incidentally the typical model for spaces with constant positive curvature (gaussian curvature $K = 1$), and it could therefore be argued that we have included examples for each distinct space of constant curvature.

BIBLIOGRAPHY

- [1] P. Arnoux. ‘Sturmian Sequences’. In: *Substitutions in Dynamics, Arithmetics and Combinatorics*. Ed. by N. Pytheas Fogg, Val  re Berth  , S  bastien Ferenczi, Christian Mauduit and Anne Siegel. Springer, Berlin, Heidelberg, 2002.
- [2] Dzimitry Badziahin, Andrew Pollington and Sanju Velani. *On a problem in simultaneous Diophantine approximation: Schmidt’s conjecture*. Read at arXiv:1001.2694 [math.NT], As of February 2020 online at <https://arxiv.org/pdf/1001.2694.pdf>.
- [3] J. W. S. Cassels. *An Introduction to Diophantine Approximation*. The Syndics of the Cambridge University Press, 1957.
- [4] Keith Conrad. *SL₂(Z)*. As of February 2020 online at [https://kconrad.math.uconn.edu/blurbs/grouptheory/SL\(2,Z\).pdf](https://kconrad.math.uconn.edu/blurbs/grouptheory/SL(2,Z).pdf).
- [5] M. Maurice Dodson and Simon Kristensen. ‘Hausdorff Dimension and Diophantine Approximation’. In: *Proceedings of Symposia In Pure Mathematics*. Read at arXiv:math/0305399 [math.NT], As of February 2020 online at <https://arxiv.org/pdf/math/0305399.pdf>.
- [6] Encyclopedia of Mathematics. *Thue–Siegel–Roth theorem*. As of February 2020 online at http://www.encyclopediaofmath.org/index.php?title=Thue-Siegel-Roth_theorem&oldid=34911.
- [7] Kenneth Falconer. *Fractal Geometry, Mathematical Foundations and Applications*. John Wiley & Sons, 1990.
- [8] Paul Garrett. *Liouville’s theorem on diophantine approximation*. As of February 2020 online at https://www-users.math.umn.edu/~garrett/m/mfms/notes_2013-14/04b_Liouville_approx.pdf. 2013.
- [9] Sigmundur Gudmundsson. *Lecture Notes in Mathematics, An Introduction to Gaussian Geometry*. version 2.0083, As of February 2020 online at <http://www.matematikk.lu.se/matematikklu/personal/sigma/Gauss.pdf>. 2019.
- [10] Oleg Karpenkov. *Geometry of Continued Fractions*. 1st ed. Springer, Berlin, Heidelberg, 2013.
- [11] OEIS Foundation Inc. (2019), The On-Line Encyclopedia of Integer Sequences. *A001203 Simple continued fraction expansion of Pi*. As of February 2020 online at <https://oeis.org/A001203>.

- [12] Jörg Schmeling, Endre Szabó and Reinhard Winkler. 'Hartman and Beatty bisequences'. In: *Algebraic Number Theory and Diophantine Analysis, Proceedings of the International Conference held in Graz, Austria August 30 to September 5, 1998*. Ed. by Franz Halter-Koch and Robert F. Tichy. Walter de Gruyter, 1998.
- [13] Caroline Series. *Hyperbolic geometry, MA 448*. Lecture notes. 2013.
- [14] Caroline Series. *Continued Fractions and Hyperbolic Geometry*. Loughborough LMS Summer School. Lecture notes. 2015.
- [15] Hanna Uscka-Wehlou. 'Digital lines, Sturmian words, and continued fractions'. As of February 2020 online at <https://www.diva-portal.org/smash/get/diva2:228516/FULLTEXT01.pdf>. PhD thesis. Uppsala University, 2009.