# CONTINUED FRACTIONS, PELL'S EQUATION, AND TRANSCENDENTAL NUMBERS

JEREMY BOOHER

Continued fractions usually get short-changed at PROMYS, but they are interesting in their own right and useful in other areas of number theory. For example, they given a way to write a prime congruent to 1 modulo 4 as a sum of two squares. They can also be used to break RSA encryption when the decryption key is too small. Our first goal will be to show that continued fractions are "the best" approximations of real numbers in a way to be made precise later. Then we will look at their connection to lines of irrational slope in the plane, Pell's Equation, and their further role in number theory.

## 1. BASIC PROPERTIES

First, let's establish notation. For $\beta \in \mathbb{R}$, let $\beta_0 := \beta$ and define

$$a_i := [\beta_i] \quad \text{and} \quad \beta_{i+1} := \frac{1}{\beta_i - a_i}.$$

The $n$th convergent to $\beta$ is the fraction

$$[a_0, a_1, \ldots, a_n] := a_0 + \cfrac{1}{a_1 + \cfrac{1}{\ldots + \cfrac{1}{a_n}}}.$$

The numerator and denominator, when this fraction is written in lowest terms, are denoted by $p_n$ and $q_n$.

By a simple induction, we have that

$$\beta = a_0 + \cfrac{1}{a_1 + \cfrac{1}{\ldots + \cfrac{1}{a_{n-1} + \frac{1}{\beta_n}}}} = [a_0, a_1, \ldots, a_{n-1}, \beta_n].$$

If at any point the remainder $\beta_i$ is an integer, this process stops. In this case we know that $\beta$ is a rational number. Likewise, it is clear that if $\beta$ is rational then this process terminates. From now on, we will usually assume that $\beta$ is irrational so that this process does not terminate.

Our first observation is about which $[a_0, a_1, \ldots, a_n]$ are less than $\beta$.

**Proposition 1.** *If $n$ is even, $[a_0, a_1, \ldots, a_n]$ is less than $\beta$, otherwise it is greater.*

*Proof.* The proof uses the following simple lemma.

**Lemma 2.** *For any $a_i$ of length $n$, if $a_n < a'_n$ and if $n$ is even then*

$$[a_0, a_1, \ldots, a_n] < [a_0, a_1, \ldots, a'_n].$$

*If $n$ is odd then the reverse inequality holds.*

*Proof.* We use induction on $n$. For $n = 0$, $a_0 < a'_0$ is the desired conclusion. Otherwise, write

$$[a_0, a_1, \ldots, a_n] = a_0 + \frac{1}{[a_1, \ldots, a_n]}.$$

---

Suppose $n$ is even. By the inductive hypothesis, the denominator is less than $[a_1, \ldots, a'_n]$. Thus

$$[a_0, a_1, \ldots, a_n] < a_0 + \frac{1}{[a_1, \ldots, a'_n]} = [a_0, a_1, \ldots, a'_n].$$

If $n$ is odd, the denominator is greater and the reverse inequality follows. $\square$

We can now prove the Proposition. Note that

$$\beta = [a_0, a_1, \ldots, \beta_n]$$

and that $\beta_n > a_n$ (we can't have equality since $\beta$ is irrational by assumption). By the lemma, if $n$ is even $[a_0, a_1, \ldots, \beta_n] > [a_0, a_1, \ldots, a_n]$ and if $n$ is odd $[a_0, a_1, \ldots, \beta_n] < [a_0, a_1, \ldots, a_n]$. $\square$

The next order of business is to derive the recursive formula for $p_n$ and $q_n$. The following method is not the most direct but will be useful later.

**Definition 3.** Define $\{a_0\} := a_0$ and $\{a_0, a_1\} := a_1 a_0 + 1$. Then inductively define

$$\{a_0, a_1, \ldots, a_n\} := \{a_0, a_1, \ldots, a_{n-1}\} a_n + \{a_0, a_1, \ldots, a_{n-2}\}.$$

The $a_i$ do not a priori need to arise from a continued fraction.

**Proposition 4** (Euler)**.** *Let $S_k$ be the set of all increasing sequences of length $n + 1 - 2k$ obtained by deleting $m$ and $m + 1$ from $\{0, 1, \ldots, n\}$ $k$ times in succession. With the convention that the empty product is $1$, define*

$$l_k(a_0, a_1, \ldots, a_n) := \sum_{(b_0, \ldots, b_{n-2k}) \in S_k} a_{b_0} a_{b_1} \ldots a_{b_{n-2k}}.$$

*Then $\{a_0, a_1, \ldots, a_n\} = \sum_{0 \le k \le n+1} l_k(a_0, a_1, \ldots, a_n)$.*

*Proof.* The proof proceeds by induction on $n$. For $n = 0$, $\{a_0\} = a_0$ and $l_0(a_0) = a_0$. For $n = 1$, $\{a_0, a_1\} = a_0 a_1 + 1$. We know that $l_0(a_0, a_1) = a_0 a_1$ and $l_1(a_0, a_1) = 1$, the empty product

In general, assume the assertion holds up to length $n$. Then $\{a_0, a_1, \ldots, a_n\}$ has length $n + 1$, but by definition it is

$$\{a_0, a_1, \ldots, a_{n-1}\} a_n + \{a_0, a_1, \ldots, a_{n-2}\}.$$

By the inductive hypothesis, this equals

$$a_n \left( \sum_{0 \le k \le n} l_k(a_0, \ldots, a_{n-1}) \right) + \left( \sum_{0 \le j \le n-1} l_j(a_0, \ldots, a_{n-2}) \right).$$

The terms in the second sum are all products of $a_i$ $(0 \le i \le n)$ where the last two terms (and possibly more) are left out, while the terms in the first sum are all products of the $a_i$ with consecutive pairs left out that including $a_n$. Every product of the $a_i$'s arising by deleting multiple consecutive terms arises through exactly one of these two ways. Thus we conclude

$$\{a_0, a_1, \ldots, a_n\} = \sum_{0 \le k \le n+1} l_k(a_0, \ldots, a_n) \qquad \square$$

*Example* 5. Although this seems complicated, all of the hardness is in the notation. For the case of $n = 3$, all this is asserting is that

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

Every term in this sum is obtained by deleting zero, one, or two pairs of consecutive $a_i$. If we look at $\{a_0, a_1, a_2, a_3, a_4\}$, when we remove pairs of consecutive terms we either remove $a_3$ and $a_4$ or we don't. If we do, then all the terms we are adding up are just terms in the sum for $\{a_0, a_1, a_2\}$. If we don't, then $a_4$ is in each of the terms and we can remove zero or more pairs of consecutive

terms from $a_0, \ldots, a_3$ and add the products up. By induction, this is $a_4\{a_0, a_1, a_2, a_3\}$. Explicitly, we have

$$\{a_0, a_1, a_2, a_3, a_4\} = a_4(a_0a_1a_2a_3 + a_0a_1 + a_0a_3 + a_2a_3 + 1) + a_0a_1a_2 + a_0 + a_2.$$

This also has a combinatorial interpretation. $\{a_0, a_1, \ldots, a_n\}$ is the number of ways to tile a 1 by $n + 1$ strip with two kinds of tiles: 1 by 2 rectangles and 1 by 1 squares, where rectangles may not overlap anything but squares may stack, with up to $a_i$ of them on the $i$th place on the strip.

**Corollary 6.** *We have that $\{a_0, a_1, a_2, \ldots, a_n\} = \{a_n, a_{n-1}, \ldots, a_0\}$.*

*Proof.* Note that the description in the Proposition depends only on which $a_i$ are consecutive, not the actual ordering. Therefore the two expressions are equal. $\square$

The next proposition will later be interpreted as a fact about determinants and about the difference between two fractions.

**Proposition 7.** *For $a_0, \ldots, a_n$, we have*

$$\{a_0, a_1, \ldots, a_{n-1}, a_n\}\{a_1, \ldots, a_{n-1}\} - \{a_1, \ldots, a_{n-1}, a_n\}\{a_0, a_1, \ldots, a_{n-1}\} = (-1)^{n+1}.$$

*Proof.* The $n = 0$ and $n = 1$ cases are trivial. By induction, suppose it holds for $n - 1$. Then using the definition of $\{a_0, \ldots, a_i\}$,

$$\begin{aligned}
&\{a_0, a_1, \ldots, a_{n-1}, a_n\}\{a_1, \ldots, a_{n-1}\} - \{a_1, \ldots, a_{n-1}, a_n\}\{a_0, a_1, \ldots, a_{n-1}\} \\
&= (a_n\{a_0, a_1, \ldots, a_{n-1}\} + \{a_0, a_1, \ldots, a_{n-2}\})\{a_1, \ldots, a_{n-1}\} \\
&\quad - (\{a_1, \ldots, a_{n-1}\}a_n + \{a_1, \ldots, a_{n-2}\})\{a_0, a_1, \ldots, a_{n-1}\} \\
&= \{a_0, a_1, \ldots, a_{n-2}\}\{a_1, \ldots, a_{n-1}\} - \{a_1, \ldots, a_{n-2}\}\{a_0, a_1, \ldots, a_{n-1}\} \\
&= -(-1)^n = (-1)^{n+1}
\end{aligned}$$

by the inductive hypothesis. $\square$

We can now prove the standard recursive formulas for $p_n$ and $q_n$.

**Proposition 8.** *Let $\beta = [a_0, a_1, a_2, \ldots]$ be a continued fraction. The numerator and denominators of the nth convergent are $\{a_0, a_1, \ldots, a_n\}$ and $\{a_1, a_2, \ldots, a_n\}$. Thus they can be calculated recursively by the formulas*

$$p_n = a_np_{n-1} + p_{n-2} \quad and \quad q_n = a_nq_{n-1} + q_{n-2}.$$

*Proof.* As usual, the proof proceeds by induction. For $n = 0$ or $n = 1$, the convergents are $a_0$ and $a_0 + \frac{1}{a_1}$, the numerators are $a_0$ and $\{a_0, a_1\} = a_0a_1 + 1$, and the denominators are 1 and $a_1$. Now assume this holds for the $n - 1$th convergent of any continued fraction. In particular, we know that

$$[a_1, a_2, \ldots, a_n] = \frac{\{a_1, a_2, \ldots, a_n\}}{\{a_2, \ldots, a_n\}}$$

since it is of length $n$ while $\dfrac{p_n}{q_n} = a_0 + \dfrac{1}{[a_1, a_2, \ldots, a_n]}$ by definition. Thus combining the fractions gives

$$\begin{aligned}
\frac{p_n}{q_n} &= \frac{a_0\{a_1, \ldots, a_n\} + \{a_2, \ldots, a_n\}}{\{a_1, \ldots, a_n\}} \\
&= \frac{a_0\{a_n, a_{n-1}, \ldots, a_1\} + \{a_n, a_{n-1}, \ldots, a_2\}}{\{a_1, \ldots, a_n\}} \qquad \text{(Corollary 6)} \\
&= \frac{\{a_n, a_{n-1}, \ldots, a_0\}}{\{a_1, \ldots, a_n\}} \qquad\qquad\qquad\qquad \text{(Definition)} \\
&= \frac{\{a_0, a_1, \ldots, a_n\}}{\{a_1, \ldots, a_n\}}. \qquad\qquad\qquad\qquad \text{(Corollary 6)}
\end{aligned}$$

To show these are in fact $p_n$ and $q_n$, we need to know they are relatively prime. This follows from Proposition 7. □

Finally, there is one more formula similar to Proposition 7 we will need.

**Proposition 9.** *For any continued fraction,*

$$p_n q_{n-2} - q_n p_{n-2} = (-1)^{n-1} a_n.$$

*Proof.* This will follow from Proposition 7. For $n \geq 2$, we calculate

$$\begin{aligned} p_n q_{n-2} - q_n p_{n-2} &= (a_n p_{n-1} + p_{n-2}) q_{n-2} - (a_n q_{n-1} + q_{n-2}) p_{n-2} \\ &= a_n (p_{n-1} q_{n-2} - p_{n-2} q_{n-1}) \\ &= (-1)^n a_n. \end{aligned}$$

□

## 2. Continued Fractions as Best Approximations

The previous algebraic work gives us plenty of information about the convergence of continued fractions.

**Theorem 10.** *The convergents $\frac{p_n}{q_n}$ to $\beta$ actually converge to $\beta$. More precisely, we know*

$$\left| \beta - \frac{p_n}{q_n} \right| < \frac{1}{q_n q_{n+1}}.$$

*Furthermore, the even convergents are less than $\beta$ and the odd convergents are greater than $\beta$.*

*Proof.* Rewriting Propositions 7 and 9 in terms of fractions, we have

$$\frac{p_n}{q_n} - \frac{p_{n-1}}{q_{n-1}} = \frac{(-1)^{n-1}}{q_n q_{n-1}} \quad \text{and} \quad \frac{p_n}{q_n} - \frac{p_{n-2}}{q_{n-2}} = \frac{(-1)^{n-1} a_n}{q_{n-2} q_n}.$$

In particular, the second shows that the sequence $\frac{p_1}{q_1}, \frac{p_3}{q_3}, \frac{p_5}{q_5}, \dots$ is a monotonic increasing sequence. Likewise, $\frac{p_0}{q_0}, \frac{p_2}{q_2}, \frac{p_4}{q_4}, \dots$ is a monotonic decreasing sequence. This implies that the sequences converge or diverge to $\pm\infty$ However, the first equation shows that the even and odd convergents become arbitrarily close, hence the two series converge to the same thing. We know that the odd convergents are greater than $\beta$ and the even ones less because of Proposition 1, so the convergents converge to $\beta$. Since consecutive convergents are on opposite sides of $\beta$,

$$\left| \frac{p_n}{q_n} - \beta \right| < \left| \frac{p_n}{q_n} - \frac{p_{n+1}}{q_{n+1}} \right| = \frac{1}{q_n q_{n+1}}.$$

□

**Corollary 11.** *With the previous notation,*

$$\left| \frac{p_n}{q_n} - \beta \right| < \frac{1}{q_n^2}.$$

*Proof.* By definition, $q_{n+1} = a_n q_n + q_{n-1} \geq 1 \cdot q_n + q_{n-1} \geq q_n$. □

*Remark* 12. Because $q_{n+1} = a_n q_n + q_{n-1} \geq q_n + q_{n-1}$, the denominators grow at least as fast as the Fibonacci numbers, so $q_n$ is exponential in $n$. Calculating just a few convergents can provide very good approximations of irrational numbers.

We can say something stronger about one of every two convergents.

**Proposition 13.** *At least one of every pair of consecutive convergents satisfies*

$$\left| \frac{p_n}{q_n} - \beta \right| < \frac{1}{2q_n^2}.$$

*Proof.* Suppose neither of $\frac{p_n}{q_n}$ and $\frac{p_{n+1}}{q_{n+1}}$ satisfy this. Then because $\beta$ lies between them we have that

$$\left|\frac{p_n}{q_n} - \frac{p_{n+1}}{q_{n+1}}\right| = \left|\frac{p_n}{q_n} - \beta\right| + \left|\frac{p_{n+1}}{q_{n+1}} - \beta\right| > 2\sqrt{\left|(\frac{p_n}{q_n} - \beta)(\frac{p_{n+1}}{q_{n+1}} - \beta)\right|}$$

by the arithmetic-geometric mean inequality. This is a strict inequality because $\frac{p_n}{q_n} - \beta$ and $\frac{p_{n+1}}{q_{n+1}} - \beta$ cannot be equal as $\beta$ is irrational. By our assumption we have

$$\left|\frac{p_n}{q_n} - \frac{p_{n+1}}{q_{n+1}}\right| = 2\sqrt{\left|(\frac{p_n}{q_n} - \beta)(\frac{p_{n+1}}{q_{n+1}} - \beta)\right|} \geq 2\sqrt{\frac{1}{2q_n^2}\frac{1}{2q_{n+1}^2}} = \frac{1}{q_n q_{n+1}}.$$

This is a contradiction with Proposition 7, which says

$$\left|\frac{p_n}{q_n} - \frac{p_{n+1}}{q_{n+1}}\right| = \frac{1}{q_n q_{n+1}}. \qquad \square$$

Infinitely many convergents also exist for which we can replace the 2 with a $\sqrt{5}$. The $\sqrt{5}$ is optimal, as can be seen by looking at $\frac{-1+\sqrt{5}}{2}$. But excluding this number, $2\sqrt{2}$ works. For more details, see Hardy and Wright.

It is also possibly so put a limit on how good an approximation a convergent can be. As always, remember that this is for irrational numbers only.

**Proposition 14.** *For any convergent $\frac{p_n}{q_n}$ to $\beta$, one has that*

$$\left|\frac{p_n}{q_n} - \beta\right| > \frac{1}{q_n(q_{n+1} + q_n)}$$

*Proof.* Since the odd and even convergents form monotonic sequences, $\frac{p_{n+2}}{q_{n+2}}$ is closer to $\beta$ than $\frac{p_n}{q_n}$ is. Thus

$$\left|\frac{p_n}{q_n} - \beta\right| > \left|\frac{p_n}{q_n} - \frac{p_{n+2}}{q_{n+2}}\right| = \frac{a_{n+2}}{q_n q_{n+2}}.$$

But $q_{n+2} = a_{n+2}q_{n+1} + q_n$, so

$$\frac{a_{n+2}}{q_{n+2}} > \frac{1}{q_{n+1} + q_n/a_{n+2}} > \frac{1}{q_{n+1} + q_n}$$

so we conclude

$$\left|\frac{p_n}{q_n} - \beta\right| > \frac{1}{q_n(q_{n+1} + q_n)}. \qquad \square$$

The next step is to investigate in what sense continued fractions are the best approximations to irrational numbers. There are several different ways to measure this. The first is simply to look at $\frac{p}{q} - \beta$ versus $\frac{1}{q^2}$, as suggested by the previous propositions.

**Theorem 15.** *Suppose $|\frac{p}{q} - \beta| < \frac{1}{2q^2}$. Then $\frac{p}{q}$ is a convergent to $\beta$.*

The proof of this will rely on a different notion of closeness that is motivated by viewing irrational numbers as slopes of lines and continued fractions as lattice points close to the line. This will be discussed in the next section. For now, we will show the following:

**Proposition 16.** *Suppose $|p - q\beta| \leq |p_n - q_n\beta|$ and $0 < q < q_{n+1}$. Then $q = q_n$ and $p = p_n$.*

*Proof.* The key fact is that the matrix

$$\begin{pmatrix} p_n & p_{n+1} \\ q_n & q_{n+1} \end{pmatrix}$$

has determinant $\pm 1$ (Proposition 7), so there are integer solutions $(u, v)$ to the system of equations $p = up_n + vp_{n+1}$ and $q = uq_n + vq_{n+1}$. Note that $uv \leq 0$: if $u$ and $v$ were of the same sign and nonzero, then $|q| > |q_{n+1}|$ which contradicts our hypothesis. Now write

$$|p - q\beta| = |u(p_n - \beta q_n) + v(p_{n+1} - \beta q_{n+1})|.$$

Since consecutive convergents lie on opposite sides of $\beta$ and $uv \leq 0$, $u(p_n - \beta q_n)$ and $v(p_{n+1} - \beta q_{n+1})$ have the same sign, or one is zero. This means

$$|p - q\beta| = |u(p_n - \beta q_n)| + |v(p_{n+1} - \beta q_{n+1})|.$$

For this to be less than $|p_n - q_n\beta|$, we must have either $|u| = 1$ and $v = 0$ or have that $u = 0$. In the latter case, $q$ is a multiple of $q_{n+1}$, a contradiction. If the former, then since $q$ is positive $u$ must be as well, so $q = uq_n + vq_{n+1} = q_n$ and $p = up_n + vp_{n+1} = p_n$ and we are done. $\square$

With this, we can prove Theorem 15.

*Proof.* Suppose $|\frac{p}{q} - \beta| < \frac{1}{2q^2}$. If $p$ and $q$ are not relatively prime, say $p = dp'$ and $q = dq'$ then

$$\left| \frac{p'}{q'} - \beta \right| < \frac{1}{2d^2(q')^2} \leq \frac{1}{2(q')^2}$$

so we may assume $p$ and $q$ are relatively prime. We may also assume that $q$ is positive by possibly changing signs. If $\frac{p}{q}$ is not a convergent, we can pick $n$ so that $q_n < q < q_{n+1}$. In the case that

$$|p - q\beta| \leq |p_n - q_n\beta|$$

then by the previous proposition $p = p_n$ and $q = q_n$. Thus we may assume that

$$|p - q\beta| \geq |p_n - q_n\beta| \quad \text{and so} \quad |p_n - q_n\beta| < \frac{1}{2q}.$$

Now we can calculate

$$\left| \frac{p}{q} - \frac{p_n}{q_n} \right| \leq \left| \frac{p}{q} - \beta \right| + \left| \frac{p_n}{q_n} - \beta \right| < \frac{1}{2q_nq} + \frac{1}{2q^2} \leq \frac{1}{2q_nq} + \frac{1}{2qq_n} = \frac{1}{qq_n}.$$

However,

$$\left| \frac{p}{q} - \frac{p_n}{q_n} \right| = \frac{|pq_n - qp_n|}{qq_n}$$

and since the numerator is either 0 or a positive integer it must be zero which implies that $\frac{p}{q}$ is a convergent to $\beta$. $\square$

This justifies the informal contention that convergents are the best approximation to irrational numbers. However, there can be other fractions which nevertheless are very good approximations as well if the meaning of "very good" is changed slightly. For example, there may be other fractions that satisfy

$$\left| \frac{p}{q} - \beta \right| < \frac{1}{2q_n^2} \quad \text{and} \quad q_n < q < q_{n+1}.$$

For example, $\frac{8}{3}$ and $\frac{37}{14}$ are consecutive convergents to $\sqrt{7}$. However, $\frac{13}{5}$ satisfies

$$\left| \frac{13}{5} - \sqrt{7} \right| < \frac{1}{2 \cdot 3^2}.$$

It turns out that $\frac{13}{5}$ arises as the term between two previous convergents $\frac{8}{3}$ and $\frac{5}{2}$ in a Farey sequence. It is their mediant and because it is a good approximation is called a semi-convergent. Note it is not as good an approximation as a convergent relative to the size of its denominator since

$$\left| \frac{13}{5} - \sqrt{7} \right| > \frac{1}{2 \cdot 5^2}.$$

## 3. Continued Fractions, Lines of Irrational Slope, and Lattice Points

A geometric way to make sense of continued fractions is to view $\beta$ as the line $y = \beta x$ passing through the origin and represent rational numbers $\frac{p}{q}$ as the lattice point $(q, p)$. The distance between a point $(x_0, y_0)$ and the line $ax + by + c = 0$ is

$$\frac{|ax_0 + by_0 + c|}{\sqrt{a^2 + b^2}}.$$

Thus the distance between $(q, p)$ and $y = \beta x$ is $\frac{|\beta q - p|}{\sqrt{\beta^2 + 1}}$. Up to a scaling factor, this is the definition of distance that appeared in Proposition 16. Reinterpreting it in the language of lines and lattice points, we see:

**Theorem 17.** *Let $y = \beta x$ be a line with irrational slope and $\frac{p_n}{q_n}$ be the nth convergent to $\beta$. If $(q, p)$ is closer to the line than $(q_n, p_n)$ and $q < q_{n+1}$, then $q = q_n$ and $p = p_n$.*

The key step in the proof, writing $p = up_n + vp_{n+1}$ and $q = uq_n + vq_{n+1}$, is just expressing the vector $(p, q)$ as a linear combination of the two nearest convergents $(p_n, q_n)$ and $(p_{n+1}, q_{n+1})$ which are generators for the standard lattice.

Interpreting other algebraic facts in this context, we see that the convergents are alternatingly above and below the line. Furthermore, the estimates on how close convergents are to $\beta$ give estimates on how well the line $px - qy = 0$ approximates the slope of $y - \beta x = 0$.

It is possible to prove all of the algebraic statements in the first section geometrically using this picture. See for example "An Introduction to Number Theory" by Harold Stark.

## 4. Pell's Equation

Continued fractions provide a way to analyze solutions to Pell's equation and its relatives $x^2 - dy^2 = r$ when $r$ is small compared to $d$. All integral solutions come from convergents to $\sqrt{d}$.

**Theorem 18.** *Let $d$ be a positive square free integer and $r \in \mathbb{Z}$ satisfy $r^2 + |r| \le d$. Suppose $x$ and $y$ are positive integers that satisfy $x^2 - dy^2 = r$. Then $\frac{x}{y}$ is a convergent to $\sqrt{d}$.*

*Proof.* First, a bit of algebra. Since $y \ge 1$, we have that $\frac{r}{y^2} + d$ is minimized when $r$ is negative and $y = 1$. Thus we have

$$\frac{\sqrt{d}}{|r|} + \frac{\sqrt{\frac{r}{y^2} + d}}{|r|} \ge \frac{\sqrt{d} + \sqrt{d - |r|}}{|r|}.$$

By hypothesis, $\sqrt{d} > |r|$ and $d - |r| \ge |r|^2$. Thus we have

$$\frac{\sqrt{d}}{|r|} + \frac{\sqrt{\frac{r}{y^2} + d}}{|r|} > 2.$$

Now since $(x - \sqrt{d}y)(x + \sqrt{d}y) = r$ we know that

$$\left| \frac{x}{y} - \sqrt{d} \right| = \frac{|r|}{|y(x + \sqrt{d}y)|} = \frac{|r|}{y^2(\sqrt{d} + \sqrt{r/y^2 + d})} < \frac{1}{2y^2}$$

by the previous algebraic computation. By Theorem 15, $\frac{x}{y}$ is a convergent to $\sqrt{d}$. $\qquad \square$

Since all small values of $x^2 - dy^2$ arise from convergents, if $x^2 - dy^2 = 1$ is to have a solution it must arise from a convergent to $\sqrt{d}$.

**Theorem 19.** *Pell's equation $x^2 - dy^2 = 1$ has a non-trivial solution for any square-free integer $d$.*

*Proof.* Any convergent $\frac{p}{q}$ to $\sqrt{d}$ satisfies

$$\left| \frac{p}{q} + \sqrt{d} \right| < 1 + 2\sqrt{d}$$

since $\frac{p}{q}$ can be at most one away from $\sqrt{d}$. Combining this with Proposition 11, we see that

$$\left| p^2 - dq^2 \right| < \frac{1 + 2\sqrt{d}}{q^2} \cdot q^2 = 1 + 2\sqrt{d}$$

There are an infinite number of convergents since $\sqrt{d}$ is irrational, and only a finite number of choices for $p^2 - dq^2$, so there must be an $r$ with an infinite number of convergents satisfying $p^2 - dq^2 = r$. There are only a finite number of choices for $(p, q)$ to reduce to modulo $r$, and an infinite number of convergents satisfying the equation, so there two distinct convergents $(p_0, q_0)$ and $(p_1, q_1)$ with

$$p_0 \equiv p_1 \mod r \quad \text{and} \quad q_0 \equiv q_1 \mod r \quad \text{and} \quad p_0^2 - dq_0^2 = p_1^2 - dq_1^2 = r.$$

Then looking at the ratio

$$u = \frac{p_0 + q_0\sqrt{d}}{p_1 + q_1\sqrt{d}} = \frac{(p_0 + q_0\sqrt{d})(p_1 - q_1\sqrt{d})}{r}$$

$$= \frac{p_0 p_1 - dq_0 q_1 + \sqrt{d}(p_1 q_0 - q_1 p_0)}{r}.$$

Since $p_0 p_1 - dq_0 q_1 \equiv p_0^2 - dq_0^2 \equiv 0 \mod r$ and $p_1 q_0 - q_1 p_0 \equiv 0 \mod r$, the ratio $u$ is of the form $p + q\sqrt{d}$ with $p, q \in \mathbb{Z}$. Since the norm from $\mathbb{Q}(\sqrt{d})$ to $\mathbb{Z}$ is multiplicative, $u$ has norm 1. Therefore Pell's equation has a non-trivial solution.[1] $\qquad\square$

Now that we have one non-trivial solution, we can determine all solutions to the equations $x^2 - dy^2 = 1$. This is easiest to understand in terms of the arithmetic of $\mathbb{Q}(\sqrt{d})$. The key fact is that the norm of $x + y\sqrt{d}$ is $x^2 - dy^2$ and is a multiplicative function $\mathbb{Q}(\sqrt{d}) \to \mathbb{Q}$. The units in the ring of integers are elements with norm $\pm 1$, so we will first investigate solutions to $x^2 - dy^2 = \pm 1$. We first define a fundamental solution.

**Definition 20.** A solution $(x, y)$ to Pell's equation $x^2 - dy^2 = \pm 1$ is a fundamental solution if $x$ and $y$ are positive, $x > 1$, and there is no solution $(x', y')$ with $x'$ and $y'$ positive and $x > x' > 1$.

A fundamental solution always exists since given any non-trivial solution we can negate $x$ and $y$ and obtain a solution with $x$ and $y$ positive; there are a finite number of smaller positive values of $x$, so there is a smallest.

**Theorem 21.** *Let $d$ be a square free integer. Let $(x, y)$ be a fundamental solution to $x^2 - dy^2 = \pm 1$. For solution $(a, b)$ to this equation there is an $n \in \mathbb{Z}$ and $\epsilon = \pm 1$ such that*

$$a + b\sqrt{d} = \epsilon(x + y\sqrt{d})^n.$$

**Corollary 22.** *Let $d \equiv 2, 3 \mod 4$ be square-free and positive, and let $K = \mathbb{Q}(\sqrt{d})$. The group of units in $\mathcal{O}_K$ is isomorphic $\mathbb{Z}/2\mathbb{Z} \times \mathbb{Z}$.*

This is Dirichlet's unit theorem for a real quadratic field.

**Corollary 23.** *Let $d$ be a square free integer. Let $(x, y)$ be a fundamental solution to $x^2 - dy^2 = 1$. For solution $(a, b)$ to this equation there is an $n \in \mathbb{Z}$ and $\epsilon = \pm 1$ such that*

$$a + b\sqrt{d} = \epsilon(x + y\sqrt{d})^n.$$

---

[1] An alternate proof will be given using Proposition 33 which is computationally helpful.

This follows from the main theorem since if the fundamental unit has norm $-1$, its square has norm 1 and generates all other solutions with norms 1.

*Proof.* To prove the theorem, suppose $a + b\sqrt{d}$ is not of the form $\pm(x + y\sqrt{d})^m$. Then pick the appropriate sign $\epsilon$ and integer $m$ so that

$$(x + y\sqrt{d})^m < \epsilon(a + b\sqrt{d}) < (x + y\sqrt{d})^{m+1}.$$

But then since

$$1 < \epsilon(a + b\sqrt{d})(x - y\sqrt{d})^m < (x + y\sqrt{d})$$

we can obtain a contradiction to $(x, y)$ being a fundamental solution. Letting $x' + y'\sqrt{d} := \epsilon(a + b\sqrt{d})(x - y\sqrt{d})^m$, since the norm is multiplicative it follows that $(x')^2 - d(y')^2 = \pm 1$. Choose signs so that $x'$ and $y'$ are positive. Then since $y^2 = \frac{x^2 \pm 1}{d}$, we have

$$1 < x' + y'\sqrt{d} = x' + \sqrt{x'^2 \pm 1} < x + \sqrt{x^2 \pm 1}.$$

Since $x'$ and $x$ are positive integers and $f(x) = x + \sqrt{x^2 \pm 1}$ is an increasing function of $x$ on this range, $x' < x$. This contradicts the assumption that $(x, y)$ was a fundamental solution. $\square$

## 5. PERIODIC CONTINUED FRACTIONS, QUADRATIC IRRATIONALS, AND THE SUPER MAGIC BOX

By a quadratic irrational, I mean a real solution $\alpha$ to a quadratic equation $ax^2 + bx + c = 0$ where $a, b, c \in \mathbb{Z}$: an element of $\mathbb{Q}(\sqrt{d})$. From experience, we know that all quadratic irrationals seems to have periodic continued fractions in the sense that there exist $N$ and $k$ such that $a_n = a_{n+k}$ for $n > N$. Conversely, all periodic continued fractions seem to arise from quadratic irrationals.

By analogy with repeating decimals, a bar indicates repeating values of $a_i$.

**Theorem 24** (Euler). *Let* $\beta = [a_0, a_1, \ldots, a_{n-1}, \overline{a_n, \ldots, a_{n+k-1}}]$ *be a periodic continued fraction. Then* $\beta$ *is a quadratic irrational.*

*Proof.* First, let $\beta_n = [\overline{a_n, \ldots, a_{n+k-1}}]$, so

$$\beta = [a_0, a_1, \ldots, a_{n-1}, \beta_n] = \frac{\beta_n p_{n-1} + p_{n-2}}{\beta_n q_{n-1} + q_{n-1}}.$$

Since $\mathbb{Q}(\sqrt{d})$ is a field, it is clear that $\beta$ is a quadratic irrational if $\beta_n$ is. But $\beta_n$ has a purely periodic continued fraction, so

$$\beta_n = [\overline{a_n, a_{n+1}, \ldots, a_{n+k-1}}] = [a_n, a_{n+1}, \ldots, a_{n+k-1}, \beta_n] = \frac{\{a_n, \ldots, a_{n+k-1}\}\beta_n + \{a_n, \ldots, a_{n+k-2}\}}{\{a_{n+1}, \ldots, a_{n+k-1}\}\beta_n + \{a_{n+1}, \ldots, a_{n+k-2}\}}.$$

In particular, we see that $\beta_n$ satisfies

$$\{a_{n+1}, \ldots, a_{n+k-1}\}\beta_n^2 + (\{a_{n+1}, \ldots, a_{n+k-2}\} - \{a_n, \ldots, a_{n+k-1}\})\beta_n - \{a_n, \ldots, a_{n+k-2}\} = 0$$

which shows $\beta_n$ and hence $\beta$ are quadratic irrationals. $\square$

Our next goal is to prove the converse.

**Theorem 25** (Lagrange). *If* $\beta > 1$ *is a quadratic irrational, then* $\beta$ *has a periodic continued fraction.*

To do this, the first step is to introduce the notion of the discriminant of a quadratic irrational, and show that all of the remainders $\beta_n$ have the same discriminant as $\beta$. Next we define what it means for a quadratic irrational to be reduced, and show there are a finite number of reduced quadratic irrationals with a specified discriminant. The last step is to show that there is an $N$ such that the remainders $\beta_n$ are reduced for $n > N$.

**Definition 26.** Let $\beta$ be a quadratic irrational which satisfies

$$A\beta^2 + B\beta + C = 0$$

where $A, B, C$ are integers with no common factors and $A > 0$. The discriminant is defined to be $B^2 - 4AC$.

**Definition 27.** For a quadratic irrational $\beta = \frac{b+\sqrt{D}}{a}$, define the conjugate to be $\beta' := \frac{b-\sqrt{D}}{a}$. A quadratic irrational $\beta$ is reduced if $\beta > 1$ and $\frac{-1}{\beta'} > 1$.

**Proposition 28.** *Let $\beta$ be a quadratic irrational with discriminant $D$. Then the remainders $\beta_n$ also have discriminant $D$.*

*Proof.* It suffices to prove that $\beta_1$ has discriminant $D$ and then use induction. Let $\beta = a_0 + \frac{1}{\beta_1}$ and let $\beta$ satisfy

$$A\beta^2 + B\beta + C = 0$$

with $A, B, C$ relatively prime integers. Then substituting and clearing the denominator shows

$$A(a_0\beta_1 + 1)^2 + B(a_0\beta_1^2 + \beta_1) + C\beta_1^2 = 0.$$

Expanding gives

$$(Aa_0^2 + Ba_0 + C)\beta_1^2 + (B + 2a_0A)\beta_1 + A = 0.$$

Note that $A$, $B + 2a_0A$, and $c + Ba_0 + Aa_0^2$ are relatively prime since $A$, $B$, and $C$ are. We may multiply by $-1$ without changing the discriminant, so we may as well assume the coefficient of $\beta_1^2$ is positive. But the discriminant of $\beta_1$ is just

$$(B + 2a_0A)^2 - 4A(Aa_0^2 + Ba_0 + C) = B^2 + 4a_0AB + 4a_0^2A^2 - 4a_0^2A^2 - 4a_0AB - 4AC = B^2 - 4AC$$

which is $D$ by definition. $\qquad\square$

**Proposition 29.** *There are only finitely many reduced quadratic irrationals of discriminant $D$.*

*Proof.* Let $\beta$ have discriminant $D$ and satisfy the polynomial $Ax^2 + Bx + C = 0$ with $A > 0$ and $A, B, C$ relatively prime integers. In other words,

$$\beta = \frac{-B + \sqrt{D}}{2A} \quad \text{and} \quad \beta' = \frac{-B - \sqrt{D}}{2A}.$$

Since it is reduced, $\beta > 1$ and $\frac{-1}{\beta'} > 1$, so we know that $0 > \beta' > -1$. In particular, this means $\beta + \beta' = \frac{-B}{C} > 0$, so $B$ must be negative. Since $\beta' < 0$, and $A > 0$,

$$\frac{-B - \sqrt{D}}{A} < 0 \Longrightarrow B > -\sqrt{D}$$

which implies there only a finite number of choices for $B$. Furthermore,

$$\beta = \frac{-B + \sqrt{D}}{2A} > 1 \Longrightarrow 2A < \sqrt{D} + \sqrt{D}$$

so there is an upper bound on $A$. But the condition

$$\beta' = \frac{-B - \sqrt{D}}{2A} > -1 \Longrightarrow -2A < -\sqrt{D}$$

give a lower bound on $A$. Thus there are finite number of choices for $A$ as well. Since $D$, $A$, and $B$ determine $C$, there are a finite number of reduced quadratic irrationals with discriminant $D$. $\quad\square$

**Proposition 30.** *For each $\beta > 1$, there exists an $N$ such that $\beta_n$ is reduced when $n > N$.*

*Proof.* All of the remainders $\beta_n$ are greater than 1, since $\beta_n = \frac{1}{\beta_{n-1}-[\beta_{n-1}]}$. The hard part is getting an expression for $\frac{-1}{\beta'_{n+1}}$ involving convergents.

Starting with the fact that

$$\beta = [a_0, \ldots, a_n, \beta_{n+1}] = \frac{p_n\beta_{n+1} + p_{n-1}}{q_n\beta_{n+1} + q_{n-1}} = \begin{pmatrix} p_n & p_{n-1} \\ q_n & q_{n-1} \end{pmatrix} \cdot \begin{pmatrix} \beta_{n+1} \\ 1 \end{pmatrix}$$

and conjugating gives that

$$\beta' = \frac{p_n\beta'_{n+1} + p_{n-1}}{q_n\beta'_{n+1} + q_{n-1}} = \begin{pmatrix} p_n & p_{n-1} \\ q_n & q_{n-1} \end{pmatrix} \cdot \begin{pmatrix} \beta'_{n+1} \\ 1 \end{pmatrix}.$$

Inverting the matrix gives that

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

and hence that

$$\frac{-1}{\beta'_{n+1}} = \frac{(p_n - q_n\beta')q_{n-1}}{(p_{n-1} - \beta'q_{n-1})q_{n-1}}.$$

The numerator can be rewritten as

$$(p_n - q_n\beta')q_{n-1} = q_n(p_{n-1} - q_{n-1}\beta') + (p_nq_{n-1} - q_np_{n-1}) = q_n(p_{n-1} - q_{n-1}\beta') + (-1)^n$$

using Proposition 7. Then

(1) $$\frac{-1}{\beta'_{n+1}} - 1 = \frac{1}{q_{n-1}}\left(q_n - q_{n-1} + \frac{(-1)^n}{(\frac{p_{n-1}}{q_{n-1}} - \beta')q_{n-1}}\right).$$

As $n$ goes to infinity in (1), $q_n - q_{n-1}$ is always a positive integer. Since $\frac{p_{n-1}}{q_{n-1}}$ approaches $\beta$ and $\beta \neq \beta'$, the denominator $(\frac{p_{n-1}}{q_{n-1}} - \beta')q_{n-1}$ goes to infinity. Thus for large $n$, the fraction is less than 1 so the right side of (1) is positive. This implies for large $n$

$$-\frac{1}{\beta'_{n+1}} - 1 > 0$$

and hence that $\beta'_{n+1}$ is reduced for large $n$. $\square$

Theorem 25 is now easy to prove.

*Proof.* If $\beta > 1$ has discriminant $D$, then there are only a finite number of reduced quadratic irrationals of discriminant $D$. All of the remainders $\beta_n$ have that discriminant, and there is an $N$ such that for $n > N$ $\beta_n$ is reduced. There are a finite number of choices for $\beta_n$, so there are $k$ and $n$ such that $\beta_n = \beta_{n+k}$. But then

$$\beta = [a_0, a_1, \ldots, a_{n-1}, \overline{a_n, \ldots, a_{n+k-1}}]$$

which shows that $\beta$ has a periodic continued fraction. $\square$

We can also say something about when continued fractions begin to be periodic.

**Theorem 31** (Galois). *Let $\beta$ be a quadratic irrational. $\beta$ is a purely periodic continued fraction if and only if $\beta$ is reduced.*

*Proof.* Suppose $\beta$ is reduced. This implies that $\beta_n$ is reduced for all $n$. By induction, it suffices to show it for $\beta_1 = \frac{1}{\beta-[\beta]}$. Since $1 > \beta - [\beta] > 0$, $\beta_1 > 1$. Since $\beta'_1 = \frac{1}{\beta'-[\beta]}$ and $0 > \beta' > -1$, $0 > \beta'_1$. Since $[\beta] \geq 1$, $\beta'_1 > -1$. Therefore all of the $\beta_n$ are reduced.

Now, suppose $\beta_{n+k} = \beta_n$. I will show that $\beta_{n+k-1} = \beta_{n-1}$ using the fact that all of the remainders are reduced. By definition,

$$\beta_n = \frac{1}{\beta_{n-1} - a_{n-1}} \quad \text{and} \quad \beta_{n+k} = \frac{1}{\beta_{n+k-1} - a_{n+k-1}}.$$

Conjugating this and rearranging gives that

$$-\frac{1}{\beta'_n} = a_{n-1} - \beta'_{n-1} \quad \text{and} \quad -\frac{1}{\beta'_{n+k}} = a_{n+k-1} - \beta'_{n+k-1}.$$

Since $\beta'_{n-1}$ and $\beta'_{n+k-1}$ are between $-1$ and $0$, $a_{n-1} = [-\frac{1}{\beta'_n}] = [-\frac{1}{\beta'_{n+k}}] = a_{n+k-1}$. Thus by induction we must have that $\beta_0 = \beta_n$ so $\beta$ is purely periodic.

Conversely, suppose a continued fraction is purely periodic. Then since the remainders are reduced for sufficiently large $n$, $\beta = \beta_k = \beta_{2k} = \dots$ is reduced. $\qquad\square$

| $A_n$ | | | 0 | 3 | 1 | 2 | 1 | 3 |
|---|---|---|---|---|---|---|---|---|
| $C_n$ | | | 1 | 4 | 3 | 3 | 4 | 1 |
| $a_n$ | | | 3 | 1 | 1 | 1 | 1 | 6 |
| $p_n$ | 0 | 1 | 3 | 4 | 7 | 11 | 18 | 119 |
| $q_n$ | 1 | 0 | 1 | 1 | 2 | 3 | 5 | 33 |
| $p_n^2 - 13q_n^2$ | | | -4 | 3 | -3 | 4 | -1 | 4 |

TABLE 1. Super Magic Box for $\sqrt{13}$

5.1. **Super Magic Box for $\sqrt{d}$.** The super magic box is an efficient computational device for computing the continued fraction of $\sqrt{d}$ for a square-free integer $d$. It is no more and no less than the standard continued fraction method with the algebra required to clear denominators replaced by simpler computational rules. Here is the standard computation to find $\sqrt{13}$'s continued fraction to compare with the super magic box. Note that $\beta_i = \frac{A_i + \sqrt{D}}{C_i}$.

$$\sqrt{3} = 3 + (\sqrt{13} - 3)$$
$$\beta_1 = \frac{3 + \sqrt{13}}{4} = 1 + \frac{\sqrt{13} - 1}{4}$$
$$\beta_2 = \frac{1 + \sqrt{13}}{3} = 1 + \frac{\sqrt{13} - 2}{3}$$
$$\beta_3 = \frac{2 + \sqrt{13}}{3} = 1 + \frac{\sqrt{13} - 1}{3}$$
$$\beta_4 = \frac{1 + \sqrt{13}}{4} = 1 + \frac{\sqrt{13} - 3}{4}$$
$$\beta_5 = \frac{3 + \sqrt{13}}{1} = 6 + \frac{\sqrt{13} - 1}{4}$$
$$\sqrt{13} = [3, \overline{1,1,1,1,6}]$$

There are many patterns visible in the super magic box. Here are few of them.

**Proposition 32.** *Let $\sqrt{d}$ have continued fraction $[a_0, \overline{a_1, \dots, a_n}]$. Then $a_n = 2a_0$.*

*Proof.* Since $\frac{1}{\sqrt{d} - [\sqrt{d}]}$ is a reduced continued fraction, it is purely periodic, so $\sqrt{d}$ is in fact of the form $[a_0, \overline{a_1, \dots, a_n}]$. Similarly, $\sqrt{d} + a_0$ has continued fraction $[2a_0, \overline{a_1, \dots, a_n}]$. However, $\sqrt{d} + a_0 > 1$ and $-1 < a_0 - \sqrt{d} < 0$ so $a_0 + \sqrt{d}$ is a reduced continued fraction and hence is purely periodic by Theorem 31. Thus $a_n = 2a_0$. $\qquad\square$

**Proposition 33.** *With the notation as in the super magic box,*

$$p_n^2 - dq_n^2 = (-1)^{n+1}C_{n+1} \quad \text{and} \quad p_n p_{n-1} - q_n q_{n-1} d = (-1)^n A_{n+1}.$$

*Proof.* This is proven by induction on $n$ for both equalities simultaneously. For $n = 0$, they can be checked directly. In general, first calculate

$$p_n p_{n+1} - q_n q_{n+1} d = a_{n+1} p_n^2 + p_{n-1} p_n - d a_{n+1} q_n^2 - d q_{n-1} q_n$$
$$= (-1)^{n+1} a_{n+1} C_{n+1} + (-1)^n A_{n+1}$$

using the inductive hypothesis. But $A_{n+2}$ is calculated in the super magic box by the formula $A_{n+2} = a_{n+1} C_{n+1} - A_{n+1}$. Thus the above is just $(-1)^{n+1} A_{n+2}$. Next calculate

$$p_{n+1}^2 - d q_{n+1}^2 = (a_{n+1} p_n + p_{n-1})^2 - d(a_{n+1} q_n + q_{n-1})^2$$
$$= a_{n+1}^2 (p_n^2 - d q_n^2) + (p_{n-1}^2 - d q_{n-1}^2) + 2a_{n+1}(p_n p_{n-1} - d q_n q_{n-1})$$
$$= (-1)^{n+1} a_{n+1}^2 C_{n+1} + (-1)^n C_n + 2a_{n+1}(-1)^n A_{n+1}$$

using the inductive hypotheses. However, by definition of $C_{n+2}$, $C_{n+1}$ and $A_{n+2}$,

$$C_{n+2} = \frac{d - A_{n+2}^2}{C_{n+1}} = \frac{d - A_{n+1}^2 + 2a_{n+1} C_{n+1} A_{n+1} - a_{n+1}^2 C_{n+1}^2}{C_{n+1}}$$
$$= \frac{C_{n+1} C_n}{C_{n+1}} + 2a_{n+1} A_{n+1} - a_{n+1}^2 C_{n+1}$$
$$= C_n + 2a_{n+1} A_{n+1} - a_{n+1}^2 C_{n+1}.$$

Substituting we get the desired result

$$p_{n+1}^2 - d q_{n+1}^2 = (-1)^n C_{n+2}. \qquad \square$$

This explains why two rows of the super magic box agree up to some signs. It also provides an alternate way to analyze Pell's equation. Suppose the continued fraction for $\beta = \sqrt{d} = [a_0, \overline{a_1, \ldots, a_n}]$. Since $\beta_n - a_n = \beta_0 - a_0$, using the notation of the super magic box we have that

$$\frac{A_n - a_n C_n + \sqrt{d}}{C_n} = \sqrt{d} - a_0.$$

Since 1 and $\sqrt{d}$ are linearly independent over the rationals, we have that $C_n = 1$. Then Proposition 33 implies that $p_{n-1}^2 - d q_{n-1}^2 = (-1)^n$. Furthermore, whenever a convergent $p_k^2 - d q_k^2 = \pm 1$, we must have that $C_k = 1$. But then $\frac{A_k + \sqrt{d}}{C_k} - a_k = \sqrt{d} - [\sqrt{d}] = \frac{1}{\beta_1}$, which shows that $k$ is a multiple of the period of $\sqrt{d}$. Thus the super magic box shows that there are solutions to Pell's equation for any $d$. The sign of the fundamental solution is determined by the parity of the length of the period for $\sqrt{d}$. Furthermore, by looking at the $p_i^2 - d q_i^2$ over the first period (if the length is even) or the first two periods (if the length is odd) and using Theorem 18 will let us find all $r$ that satisfy $r^2 + |r| \le d$ for which $x^2 - dy^2 = r$ has a non-trivial integral solution.

## 6. Other Applications of Continued Fractions

Continued fractions crop up in many areas of number theory besides the standard application to Pell's equation. They can be used to break RSA encryption if the decryption key is too small and to prove the two squares theorem.

6.1. **RSA Encryption.** In RSA encryption, Bob picks two large primes $p$ and $q$ that satisfy $p < q < 2p$, and let $n = pq$. This should be the case when doing cryptography, since there are specialized factoring algorithms that can exploit when $n$ is a product or primes of significantly different magnitude. Bob picks encryption and decryption keys $e$ and $d$ that satisfy $e = d^{-1}$ mod $\phi(n)$ using his factorization of $n$. Bob publishes $n$ and $e$, but keeps $d$, $p$, and $q$ secret. To encrypt a message $M$, Alice encodes it as a number modulo $n$ and gives Bob $C = M^e \mod n$. Bob calculates $C^d = M^{ed} = M \mod n$ (Euler's theorem) to decrypt the message. There is no known way to recover the message in general without factoring $n$, and no known way to factor $n$ efficiently.

However, if by chance $3d < n^{\frac{1}{4}}$, an adversary using knowledge of $e$ and $n$ can find $d$. Let $k = \frac{ed-1}{\phi(n)}$. Since $e < \phi(n)$, $k < d$. Since $q < \sqrt{pq}$ and $p < \sqrt{2pq}$ by hypothesis, $p + q < 3\sqrt{n}$. Thus

$$\left| \frac{e}{n} - \frac{k}{d} \right| \le \frac{|k\phi(n) + 1 - nk|}{nd} = \frac{k(p+q-1)+1}{nd} \le \frac{3k}{d\sqrt{n}} < \frac{1}{3d^2}.$$

This inequality implies that $\frac{k}{d}$ is a convergent to $\frac{e}{n}$ by Theorem 15. Using the publicly available $e$ and $n$, an adversary can use the Euclidean algorithm to find all the convergents with denominator less than $n$ in time poly-logarithmic in $n$. For each convergent, use its numerator and denominator as a guess for $k$ and $d$, and calculate what $\phi(n)$ should be. Since $p$ and $q$ satisfy the quadratic $x^2 - (n - \phi(n) + 1)x + n$, the correct guess of $\phi(n)$ will give the factorization of $n$.

### 6.2. **Sums of Two Squares.**
Another classical question in number theory is which positive primes are sums of two squares. It is easy to see by reducing modulo 4 that if $p \equiv 3 \mod 4$ it cannot be the sum of two squares. Using continued fractions, we can show that $p \equiv 1 \mod 4$ then $p$ can be written as a sum of two squares.

The idea is to look at fractions of the form $\frac{p}{q}$ where $2 \le q \le \frac{p-1}{2}$. There are two ways to write the continued fraction for a rational number: $[a_0, \ldots, a_n]$ and $[a_0, \ldots, a_{n-1}, a_n - 1, 1]$. Always use the one with $a_n \ne 1$, which is the one that comes from the Euclidean algorithm. Let the continued fraction of $\frac{p}{q}$ be $[a_0, a_1, \ldots, a_n]$. Note that $a_0 \ge 2$ since $\frac{p}{q} \ge 2$, and by our convention $a_n \ge 2$. Now consider the continued fraction $[a_n, a_{n-1}, \ldots, a_0]$. It has numerator $p$ since $\{a_n, a_{n-1}, \ldots, a_0\} = \{a_0, a_1, \ldots, a_n\} = p$. Its denominator is an integer $q' = \{a_{n-1}, \ldots, a_0\}$. Note that because $a_n \ge 2$, $\frac{p}{q'} \ge 2$ so $q' < \frac{p-1}{2}$. Also, $q \ne 1$ since $a_0 \ne 1$. Thus $[a_n, \ldots, a_0] = \frac{p}{q'}$ is another fraction of the same form as $\frac{p}{q}$. Obviously if we reverse the continued fraction of $\frac{p}{q'}$ we end up back at $\frac{p}{q}$.

Since $p \equiv 1 \mod 4$, there are $\frac{p-1}{2} - 1$ such fractions, an odd number. Since they are paired up by reversing the fraction, there must be a $q$ such that $\frac{p}{q} = [a_0, a_1, \ldots, a_{n-1}, a_n] = [a_n, a_{n-1}, \ldots, a_1, a_0]$ so $a_i = a_{n-i}$ for all $0 \le i \le n$. Now by Proposition 7,

$$p \cdot \{a_1, \ldots, a_{n-1}\} + \{a_1, \ldots, a_n\}\{a_0, a_1, \ldots, a_{n-1}\} = (-1)^n \implies p \mid \{a_0, a_1, \ldots, a_{n-1}\}^2 + (-1)^{n-1}$$

and thus $n$ is odd because $x^2 + (-1)^{n-1} = 1 \mod 4$ if and only if $(-1)^{n-1} = 1$.

Next, note that for $0 \le m < n$ we know

$$\{a_0, \ldots, a_n\} = \{a_0, a_1, \ldots, a_m\}\{a_{m+1}, \ldots, a_n\} + \{a_0, \ldots, a_{m-1}\}\{a_{m+2}, \ldots, a_n\}$$

by Proposition 4 (the first term is the terms of the sum that don't remove the pair $a_m, a_{m+1}$, the second are those terms that do). In our case, since $n$ is odd we can take $m = \frac{n-1}{2}$ and exploit the symmetry, getting

$$p = \{a_0, \ldots, a_n\} = \{a_0, \ldots, a_{(n-1)/2}\}^2 + \{a_0, \ldots, a_{(n-3)/2}\}^2.$$

Thus if $p \equiv 1 \mod 4$, $p$ is a sum of two squares.

Although this seems to give an explicit formula for the squares, it is not computationally very nice. To use it, we would first need to search through the continued fractions of all fractions of the form $\frac{p}{q}$ with $2 \le q \le \frac{p-1}{2}$ until we found a symmetric one, which without further information would be computationally expensive. However, note that the denominator $q = \{a_0, a_1, \ldots, a_{n-1}\}$ satisfies $\{a_0, a_1, \ldots, a_{n-1}\}^2 \equiv -1 \mod p$. If we could efficiently calculate a square root of $-1$ modulo $p$, we could find the two possible values for $q$ and then use the Euclidean algorithm to compute the appropriate convergents. There are general ways to do this (for example the Tonelli-Shanks algorithm), but for $-1$ it is very easy. Given a quadratic non-residue $a$, chosen by checking random integers using quadratic reciprocity, Euler's criteria says that $a^{\frac{p-1}{2}} \equiv -1 \mod p$ so $a^{\frac{p-1}{4}}$ is a square root of $-1$. Calculating $\pm a^{\frac{p-1}{4}} \mod p$ by repeated squaring gives the two possible values of $q$.

6.3. **Another Approach to the Sum of 2 Squares.** There is another approach to proving that primes congruent to one modulo four are a sum of two square using the fact that continued fractions are good approximations. Here is the key lemma, which is a slight restatement of Theorem 10.

**Lemma 34.** *For any $\beta$, not necessarily irrational, and any positive integer $n$, there exists a fraction $\frac{a}{b}$ in lowest terms so $0 < b \leq n$ and*

$$\left|\beta - \frac{a}{b}\right| < \frac{1}{b(n+1)}.$$

*Proof.* If $\beta$ is irrational, there are infinitely many convergents so we may pick $m$ so that the convergent $\frac{p_m}{q_m}$ satisfies $q_m \leq n < q_{m+1}$. Then by Theorem 10

$$\left|\beta - \frac{p_m}{q_m}\right| < \frac{1}{q_m q_{m+1}} \leq \frac{1}{q_m(n+1)}.$$

If $\beta$ is rational, the above will work unless $n$ is greater than all of the $q_m$. In that case, let $\frac{a}{b} = \beta$.  $\square$

Now, it $p$ is a prime congruent to one modulo four, then there is an integer $y$ such that $y^2 = -1$ mod $p$. Let $n = [\sqrt{p}]$. Pick a fraction $\frac{a}{b}$ with $b \leq [\sqrt{p}]$ so that

$$\left|-\frac{y}{p} - \frac{a}{b}\right| < \frac{1}{([\sqrt{p}]+1)b} < \frac{1}{\sqrt{p}b}.$$

But we also have that

$$\left|\frac{y}{p} + \frac{a}{b}\right| = \left|\frac{yb+ap}{bp}\right| = \frac{1}{\sqrt{p}b}\frac{|yb+ap|}{\sqrt{p}} < \frac{1}{\sqrt{p}b}$$

which implies that $c := yb + ap$ satisfies $|c| < \sqrt{p}$. But then

$$0 < b^2 + c^2 < 2p \quad \text{and} \quad b^2 + c^2 \equiv b^2 + y^2 b^2 \equiv 0 \mod p$$

which implies $b^2 + c^2 = p$. Thus $p$ is a sum of two squares.

6.4. **Recognizing Rational Numbers.** Continued fractions also give a way to recognize decimal approximations of rational numbers. Since a rational number has a finite continued fraction, to check whether a given decimal approximation probably comes from a rational number, run the continued fraction algorithm on the decimal approximation. If the decimal is approximating a rational, when the continued fraction algorithm should have terminated after the $n$th step, there will instead be a very tiny error between $[a_0, a_1, \ldots, a_n]$ and the decimal approximation. This will result in a huge value for the $a_{n+1}$. Looking for huge $a_i$ provides a way to find possible rational numbers that the decimal would be approximating.

For example, a simple calculation shows that $\dfrac{1003}{957} = [1, 20, 1, 4, 9]$. Approximating the fraction to 100 binary digits gives

$$1.0480668756530825496342737722$$

Changing the last digit to a 3 and running the continued fraction algorithm (with a computer, of course) gives $[1, 20, 1, 4, 9, 1078999383803443747479169]$, so we can identify it as the fraction $\frac{1003}{957}$. It is amusing that the fraction is identified although the decimal expansion has not started repeating.