

PELL'S EQUATION

1. INTRODUCTION

Fix an integer d . Our task is to find all integer solutions (x, y) to the equation

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

1.1. History. Leonhard Euler dubbed the equation **Pell's Equation** after the English mathematician John Pell (1611-1685). This terminology has persisted to the present day, despite the fact that it is "well known" to be mistaken: Pell's only contribution to the subject was the publication of some (presumably partial) results of Wallis and Brouncker. In fact the correct names are the usual ones: the problem of solving the equation was first considered by Fermat, and a complete solution was given by Lagrange.

By any name, the equation is an important one for several reasons – only some of which will be touched upon here – and its solution furnishes an ideal introduction to a whole branch of number theory, **Diophantine Approximation**.

1.2. First remarks on Pell's equation. We always have the two solutions $(x, y) = (\pm 1, 0)$, which we shall call "trivial." Recall that even having one solution is enough for us to know that there are infinitely many *rational solutions* $(x, y) \in \mathbb{Q}^2$ and indeed to find all such: we draw all lines through a single point, say $(-1, 0)$, with rational slope r and calculate the second intersection point (x_r, y_r) of this line with the quadratic equation (1). In so doing we generate a set of solutions which certainly contains all integer solutions, but figuring out which of the rational solutions are integral is by no means easy. This is a case where the question of integral solutions is essentially different, and more interesting, than the question of rational solutions.

The case $d = 0$ is trivial: clearly $(\pm 1, y)$ are the solutions.

The equation can only have solutions (x, y) with $x = 0$ if $d = -1$. Indeed, we have already made the easy observation that the solutions to $x^2 + y^2 = 1$ are $(\pm 1, 0)$ and $(0, \pm 1)$. Recall that this can be expressed by saying that the units in the ring $\mathbb{Z}[i]$ of Gaussian integers are ± 1 and $\pm i$.

So in all other cases we are reduced to looking for solutions with $xy \neq 0$. Clearly such solutions come in quadruples: if (x, y) is any such, so is $(-x, y)$, $(x, -y)$ and $(-x, -y)$. Let us therefore agree to concentrate on solutions with x and y both *positive* integers: "positive solutions."

If $d < -1$ then there are no positive solutions: $x^2 - dy^2 \geq -d > 1$. So we may assume d is a positive integer.

If $d = N^2$, then $x^2 - dy^2$ factors as $(x + Ny)(x - Ny) = 1$, and this necessitates either:

$$x + Ny = x - Ny = 1$$

in which case $x = 1, y = 0$; or

$$x + Ny = x - Ny = -1,$$

in which case $x = -1, y = 0$: there are no nontrivial solutions.

So we shall impose the condition that d be a positive, nonsquare, integer.

2. EXAMPLE: THE EQUATION $x^2 - 2y^2 = 1$

For instance, let us take $d = 2$. The equation $x^2 - 2y^2 = 1$ can be rewritten as

$$y^2 = \frac{x^2 - 1}{2};$$

in other words, we are looking for positive integers x for which $\frac{x^2-1}{2}$ is an integer square. First of all $x^2 - 1$ must be even, so x must be odd. Trying $x = 1$ gives, of course, the trivial solution $(1, 0)$. Trying $x = 3$ we get

$$\frac{3^2 - 1}{2} = 4 = 2^2,$$

so $(3, 2)$ is a nontrivial solution (our first!). Trying successively $x = 5, 7, 9$ and so forth we see that it is quite rare for $\frac{x^2-1}{2}$ to be a square: the first few values are 12, 24, 40, 60, 69, 112 and then finally with $x = 17$ we are in luck:

$$\frac{17^2 - 1}{2} = 144 = 12^2,$$

so $(17, 12)$ is another positive solution. At this point we may judge the task to be more worthy of a computer than ourselves. Indeed, our computer has no trouble finding some more solutions, which however appear to be growing quite rapidly: the next few are $(99, 70)$, $(577, 408)$, and $(3363, 2378)$.

2.1. Return of abstract algebra. If we've been paying attention, there are some clues here that we should be considering things from an algebraic perspective. Namely, (i) we see that factorization of the left-hand side of the equation $x^2 - dy^2 = 1$ leads only to trivial solutions; and (ii) when $d < 0$, we reduce to a problem that we have already solved: namely, finding all the *units* in the quadratic ring

$$\mathbb{Z}[\sqrt{d}] = \{a + b\sqrt{d} \mid a, b \in \mathbb{Z}\}.$$

We brought up the problem of determining the units in real quadratic rings $\mathbb{Z}[\sqrt{d}]$, but we did not solve it. We are coming to grips with this same problem here.

Namely, let $\alpha = r + s\sqrt{d} \in \mathbb{Q}(\sqrt{d})$ (that is, we allow r and s to be rational numbers). We put

$$\alpha' = r - s\sqrt{d},$$

and refer to α' as the **conjugate** of α . We may view conjugation as giving a homomorphism of fields from $\mathbb{Q}(\sqrt{d})$ to itself, that is an *automorphism* of $\mathbb{Q}(\sqrt{d})$. A consequence of this is that if $P(t)$ is any polynomial with rational coefficients

and $\alpha \in \mathbb{Q}(\sqrt{d})$, then $P(\alpha') = P(\alpha)$.¹

We also define a **norm map**

$$N : \mathbb{Q}(\sqrt{d}) \rightarrow \mathbb{Q}, N(\alpha) = \alpha\alpha';$$

more explicitly, $N(r + s\sqrt{d}) = (r + s\sqrt{d})(r - s\sqrt{d}) = r^2 - ds^2$. From this description it is immediate that the norm maps elements of $\mathbb{Z}[\sqrt{d}]$ to integers.

Lemma 1. *The norm map is multiplicative: for any $\alpha, \beta \in \mathbb{Q}(\sqrt{d})$, we have*

$$N(\alpha\beta) = N(\alpha)N(\beta).$$

Proof: This is a straightforward, and familiar, computation. Alternately, we get a more conceptual proof by using the homomorphism property of conjugation (which is itself verified by a simple computation!):

$$N(\alpha\beta) = \alpha\beta(\alpha\beta)' = \alpha\beta\alpha'\beta' = (\alpha\alpha')(\beta\beta') = N(\alpha)N(\beta).$$

Thus an important observation is that integral solutions (x, y) to the Pell equation (1) correspond to elements $x + y\sqrt{d}$ of $\mathbb{Z}[\sqrt{d}]$ of norm one:

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

Moreover, every norm one element $\alpha \in \mathbb{Z}[\sqrt{d}]$ is a unit: the equation $\alpha\alpha' = 1$ shows that the inverse of α is the integral element α' . We remark in passing that depending on d there may, or may not, be units of norm -1 : that is, the equation $x^2 - dy^2 = -1$ is significantly trickier to solve.

2.2. The solution for $d = 2$. An important property of the units of any ring is that they form a group under multiplication. In our context this means that if (x_1, y_1) and (x_2, y_2) are two solutions to $x^2 - dy^2 = 1$, then we can get another solution by multiplication in $\mathbb{Z}[\sqrt{d}]$:

$$N((x_1 + y_1\sqrt{d})(x_2 + y_2\sqrt{d})) = N(x_1 + y_1\sqrt{d})N(x_2 + y_2\sqrt{d}) = (x_1^2 - dy_1^2)(x_2^2 - dy_2^2) = 1 \cdot 1 = 1;$$

multiplying out $(x_1 + y_1\sqrt{d})(x_2 + y_2\sqrt{d})$ and collecting rational and irrational parts, we get a new solution $(x_1x_2 + dy_1y_2, x_1y_2 + x_2y_1)$.

Let us try out this formula in the case $d = 2$ for $(x_1, y_1) = (x_2, y_2) = (3, 2)$. Our new solution is $(3 \cdot 3 + 2 \cdot 2 \cdot 2, 2 \cdot 3 + 3 \cdot 2) = (17, 12)$, nothing else than the second smallest positive solution! If we now apply the formula with $(x_1, y_1) = (17, 12)$ and $(x_2, y_2) = (3, 2)$, we get the next smallest solution $(99, 70)$.

Indeed, for any positive integer n , we may write the n th power $(3 + 2\sqrt{2})^n$ as $x_n + y_n\sqrt{2}$ and know that (x_n, y_n) is a solution to the Pell equation. One can see from the formula for the product that it is a positive solution. Moreover, the solutions are all different because the real numbers $(3 + 2\sqrt{2})^n$ are all distinct: the only complex numbers z for which $z^n = z^m$ for some $m < n$ are the roots of unity, and the only real roots of unity are ± 1 . Indeed, we get the trivial solution $(1, 0)$ by taking the 0th power of $3 + 2\sqrt{2}$. Moreover, $(3 + 2\sqrt{2})^{-1} = 3 - 2\sqrt{2}$ is a “half-positive” solution, and taking negative integral powers of $3 + 2\sqrt{2}$ we get infinitely many more such solutions.

¹This same idea justifies the mantra “Complex roots (of real polynomials) occur in conjugate pairs” from high school algebra.

In total, every solution to $x^2 - dy^2 = 1$ that we have found is of the form $\pm(x_n, y_n)$ where $x_n + y_n\sqrt{d} = (3 + 2\sqrt{2})^n$ for some $n \in \mathbb{Z}$.

Let us try to prove that these are *all* the integral solutions. It is enough to show that every positive solution is of the form (x_n, y_n) for some positive integer n , since every norm one element $x + y\sqrt{d}$ is obtained from an element with $x, y \in \mathbb{Z}^+$ by multiplying by -1 and/or taking the reciprocal.

Lemma 2. *Let (x, y) be a nontrivial integral solution to $x^2 - dy^2 = 1$.*

- a) *x and y are both positive $\iff x + y\sqrt{d} > 1$.*
- b) *$x > 0$ and $y < 0 \iff 0 < x + y\sqrt{d} < 1$.*
- c) *$x < 0$ and $y > 0 \iff -1 < x + y\sqrt{d} < 0$.*
- d) *x and y are both negative $\iff x + y\sqrt{d} < -1$.*

Proof: Exercise.

We now observe that $3 + 2\sqrt{2}$ is the smallest positive integral solution. Note that for any positive integers x, y , $x + y\sqrt{d} \geq 1 + \sqrt{d}$. In fact, as x and y range over all positive integers, we may view $x + y\sqrt{d}$ as a double sequence of real numbers, and this double sequence tends to ∞ in the sense that for any real number M there are only finitely many terms $x + y\sqrt{d} \leq M$: indeed, since the inequality implies $x, y \leq M$ there are at most M^2 solutions. In the present case, one checks that among the elements $x + y\sqrt{d} < 3 + 2\sqrt{2} < 6$ with $x, y \in \mathbb{Z}^+$, only $2 + 3\sqrt{2}$ satisfies $x^2 - dy^2 = 1$, so it is indeed the smallest positive solution.

Now suppose we had a positive solution (x, y) which was not of the form $(3 + 2\sqrt{2})^n$. Choose the largest $n \in \mathbb{N}$ such that $x + y\sqrt{d} > (3 + 2\sqrt{2})^n$; then

$$\alpha = (x + y\sqrt{d}) \cdot (3 + 2\sqrt{2})^{-n} = (x + y\sqrt{d}) \cdot (3 - 2\sqrt{2})^n.$$

Write $\alpha = x' + y'\sqrt{d}$; then by multiplicativity (x', y') is an integral solution of the Pell equation. Moreover, by construction we have $\alpha > 1$, so by the Lemma this means that (x', y') is a positive solution. On the other hand we have that $\alpha < 3 + 2\sqrt{2}$, since otherwise we would have $x + y\sqrt{2} > (3 + 2\sqrt{2})^{n+1}$. But this is a contradiction, since we noted above that $3 + 2\sqrt{2}$ is the smallest solution with $x, y > 0$. This completes the proof.²

Thus we have “solved the Pell equation” for $d = 2$. To add icing, we can give explicit formulas for the solutions. Namely, we know that every positive integral solution (x, y) is of the form

$$x_n + y_n\sqrt{d} = (3 + 2\sqrt{2})^n$$

for $n \in \mathbb{Z}^+$. If we apply conjugation to this equation, then using the fact that it is a field homomorphism, we get

$$x_n - y_n\sqrt{d} = (3 - 2\sqrt{2})^n.$$

Let us put $u = 3 + 2\sqrt{2}$ and $u' = u^{-1} = 3 - 2\sqrt{2}$. Then, adding the two equations and dividing by 2 we get

$$x_n = \frac{1}{2}(u^n + (u')^n),$$

²We remark that this argument is reminiscent of the proof that \mathbb{Z} is a principal ideal domain.

and similarly we can solve for y_n to get

$$y_n = \frac{1}{2\sqrt{d}}(u^n - (u')^n).$$

In fact, we can do even better: $u' = 0.17157\dots$, so that for all $n \in \mathbb{Z}^+$, $\frac{1}{2}(u')^n$ and $\frac{1}{2\sqrt{d}}(u')^n$ are less than $\frac{1}{2}$. Since x_n and y_n are integers (even though the formula does not make this apparent!), this means that we can neglect the $(u')^n$ term entirely and just round to the nearest integer. For a real number α such that $\alpha - \frac{1}{2}$ is not an integer, there is a unique integer nearest to α which we shall denote by $\langle \alpha \rangle$. We have proved:

Theorem 3. *Every positive integer solution to $x^2 - 2y^2 = 1$ is of the form*

$$x_n = \left\langle \frac{(3 + 2\sqrt{2})^n}{2} \right\rangle,$$

$$y_n = \left\langle \frac{(3 + 2\sqrt{2})^n}{2\sqrt{2}} \right\rangle$$

for some positive integer n .

Among other things, this explains why it was not so easy to find solutions by hand: the size of both the x and y coordinates grow exponentially! The reader is invited to plug in a value of n for herself: for e.g. $n = 17$ it is remarkable how close the irrational numbers $u^{17}/2$ and $u^{17}/(2\sqrt{2})$ are to integers:

$$u^{17}/2 = 5168247530882.9999999999999949;$$

$$u^{17}/(2\sqrt{2}) = 3654502875938.0000000000000032.$$

A bit of reflection reveals that this has a lot to do with the fact that $\frac{x_n}{y_n}$ is necessarily very close to $\sqrt{2}$. Indeed, by turning this observation on its head we shall solve the Pell equation for general nonsquare d .

3. A RESULT OF DIRICHLET

Although we admit it is not yet clear why, in order to find all solutions of Pell's equation for general d we will need the following simple and important result.

Lemma 4. (*Dirichlet*) *For any irrational (real) number α , there are infinitely many rational numbers $\frac{x}{y}$ (with $\gcd(x, y) = 1$) such that*

$$|x/y - \alpha| < \frac{1}{y^2}.$$

Proof: Since the lowest-term denominator of any rational number $\frac{x}{y}$ is unchanged by subtracting any integer n , by subtracting the integer part $[\alpha]$ of α we may assume $\alpha \in [0, 1)$. Now divide the half-open interval $[0, 1)$ into n equal pieces: $[0, \frac{1}{n}) \cup [\frac{1}{n}, \frac{2}{n}) \dots \cup [\frac{n-1}{n}, 1)$.

Now consider the fractional parts of $0, \alpha, 2\alpha, \dots, n\alpha$. Since we have $n+1$ numbers in $[0, 1)$ and only n subintervals, by the pigeonhole principle some two of them must lie in the same subinterval. That is, there exist $0 \leq j < k \leq n$ such that

$$|j\alpha - [j\alpha] - (k\alpha - [k\alpha])| < \frac{1}{n}.$$

Now take $y = j - k$, $x = [k\alpha] - [j\alpha]$, so that the previous inequality becomes

$$|x - y\alpha| < \frac{1}{n}.$$

We may assume that $\gcd(x, y) = 1$, since were there a common factor, we could divide through by it and that would only improve the inequality. Moreover, since $0 < y < n$, we have

$$\left| \frac{x}{y} - \alpha \right| < \frac{1}{ny} < \frac{1}{y^2}.$$

This exhibits one solution. To see that there are infinitely many, observe that since α is irrational, $|\frac{x}{y} - \alpha|$ is always strictly greater than 0. But by choosing n sufficiently large we can apply the argument to find a rational number $\frac{x'}{y'}$ such that

$$\left| \frac{x'}{y'} - \alpha \right| < \left| \frac{x}{y} - \alpha \right|,$$

and hence there are infinitely many.

Remark: The preceding argument is perhaps the single most famous application of the pigeonhole principle. Indeed, in certain circles, the pigeonhole principle goes by the name “Dirichlet’s box principle” because of its use in this argument.

4. EXISTENCE OF NONTRIVIAL SOLUTIONS

Now let us prove the following result.

Theorem 5. *For any positive nonsquare integer d , the equation $x^2 - dy^2 = 1$ has a nontrivial integral solution (x, y) .*

Interestingly, the first step is to prove an “approximation”:

Proposition 6. *For some real number M , there exist infinitely many pairs of coprime positive integers (x, y) such that $|x^2 - dy^2| < M$.*

Proof: We will apply the Lemma of the preceding section to $\alpha = \sqrt{d}$: we get an infinite sequence of coprime positive (since \sqrt{d} is positive) integers (x, y) with $|\frac{x}{y} - \sqrt{d}| < \frac{1}{y^2}$. Multiplying through by y , the inequality is equivalent to

$$|x - y\sqrt{d}| < \frac{1}{y}.$$

Since

$$|x^2 - dy^2| = |x - y\sqrt{d}||x + y\sqrt{d}|,$$

in order to bound the left-hand side we also need a bound on $|x + y\sqrt{d}|$. There is no reason to expect that it is especially small, but using the triangle inequality we can get the following:

$$|x + \sqrt{d}y| = |x - \sqrt{d}y + 2\sqrt{d}y| \leq |x - \sqrt{d}y| + 2\sqrt{d}y < \frac{1}{y} + 2\sqrt{d}y.$$

Thus

$$|x^2 - dy^2| < \left(\frac{1}{y}\right)\left(\frac{1}{y} + 2\sqrt{d}y\right) = \frac{1}{y^2} + 2\sqrt{d} \leq 1 + 2\sqrt{d} = M.$$

Now let us prove Theorem 5, i.e., show that there is a nontrivial solution to $x^2 - dy^2 = 1$. We begin by further exploiting the pigeonhole principle. Namely, since we have infinitely many solutions (x, y) to $|x^2 - dy^2| < M$, there must exist some integer m ,

$|m| < M$ for which we have infinitely many solutions to the equality $x^2 - dy^2 = m$. And once again: we must have two different solutions, say (X_1, Y_1) and (X_2, Y_2) with $X_1 \equiv X_2 \pmod{|m|}$ and $Y_1 \equiv Y_2 \pmod{|m|}$ (since there are only m^2 different options altogether for $(x \pmod{|m|}, y \pmod{|m|})$ and infinitely many solutions). Let us write

$$\alpha = X_1 + Y_1\sqrt{d}$$

and

$$\beta = X_2 + Y_2\sqrt{d};$$

we have $N(\alpha) = N(\beta) = m$. A first thought is to divide α by β to get an element of norm 1; however, $\alpha/\beta \in \mathbb{Q}(\sqrt{d})$ but does not necessarily have integral x and y coordinates. However, it works after a small trick: consider instead

$$\alpha\beta' = X + Y\sqrt{d}.$$

I claim that both X and Y are divisible by m . Indeed we just calculate, keeping in mind that modulo m we can replace X_2 with X_1 and Y_2 with Y_1 :

$$X = X_1X_2 - dY_1Y_2 \equiv X_1^2 - dY_1^2 \equiv 0 \pmod{|m|},$$

$$Y = X_1Y_2 - X_2Y_1 \equiv X_1Y_1 - X_1Y_1 \equiv 0 \pmod{|m|}.$$

Thus $\alpha\beta' = m(x + y\sqrt{d})$ with $x, y \in \mathbb{Z}$. Taking norms we get

$$m^2 = N(\alpha)N(\beta') = N(\alpha\beta') = N(m(x + y\sqrt{d})) = m^2(x^2 - dy^2).$$

Since $m \neq 0$ (why?), this gives

$$x^2 - dy^2 = 1.$$

Moreover $y \neq 0$: if $y = 0$ then the irrational part of Y , namely $X_1Y_2 - X_2Y_1$, would be zero, i.e., $\frac{X_1}{Y_1} = \frac{X_2}{Y_2}$, but this is impossible since $(X_1, Y_1) \neq (X_2, Y_2)$ are both coprime pairs: they cannot define the same rational number. We are done.

5. THE MAIN THEOREM

Finally we are ready to state our main result, which determines all integral solutions to the Pell Equation.

Theorem 7. *Let d be a positive, nonsquare integer.*

- a) *There exists a positive integral solution (x, y) to $x^2 - dy^2 = 1$.*
- b) *There exists a unique positive integral solution (x_1, y_1) with $x_1 + y_1\sqrt{d}$ minimal. Put $u = x_1 + y_1\sqrt{d}$. Then every positive integral solution is of the form*

$$\left(\frac{u^n + (u')^n}{2}, \frac{u^n - (u')^n}{2\sqrt{d}} \right) = \left(\left\langle \frac{u^n}{2} \right\rangle, \left\langle \frac{u^n}{2\sqrt{d}} \right\rangle \right)$$

for a unique $n \in \mathbb{Z}^+$.

- c) *Every solution to the Pell equation is of the form $\pm(x_n, y_n)$ for $n \in \mathbb{Z}$.*

Proof: In the previous section we showed the existence of a positive solution (x, y) . It is easy to see that for any $M > 0$ there are only finitely many pairs of positive integers such that $x + y\sqrt{d} \leq M$, so among all positive solutions, there must exist one with $x + y\sqrt{d}$ least. By taking positive integral powers of this fundamental solution $x_1 + y_1\sqrt{d}$ we get infinitely many positive solutions, whose x and y coordinates can be found explicitly as in §2. Moreover, the argument of §2 – given there for $d = 2$ – works generally to show that every positive solution is of this form. The reader is invited to look back over the details.

6. A CAVEAT

We have just seen a beautiful theorem and a beautiful proof. However, it is time to admit that the phrase “solving the Pell equation” is generally taken to mean explicitly finding the fundamental solution $x_1 + y_1\sqrt{d}$. As usual in this course, we have concentrated on existence and not considered the question of how difficult it would be in practice to find the solution. Well, of course knowing that it exists we can just look for it by trying all pairs (x, y) in order of increasing $x + y\sqrt{d}$ until we find one. When $d = 2$ this was immediate. If we try other values of d we will see that sometimes it is no trouble at all:

For $d = 3$, the fundamental solution is $(2, 1)$. For $d = 6$, it is $(5, 2)$. Similarly the fundamental solution can be found by hand for $d \leq 12$; it is no worse than $(19, 6)$ for $d = 10$. However, for $d = 13$ it is $(649, 180)$: a big jump!

If we continue to search we find that the size of the fundamental solution seems to obey no reasonable law: it does not grow in a steady way with d – e.g. for $d = 42$ it is the tiny $(13, 2)$ – but sometimes it is very large: for $d = 46$ it is $(24335, 3588)$, and – hold on to your hat! – for $d = 61$ the fundamental solution is

$$(1766319049, 226153980).$$

And things get worse from here on in: one cannot count on a brute-force search for d even of modest size (e.g. five digits).

There are known algorithms which find the fundamental solution relatively efficiently. The most famous and elementary of them is as follows: one can find the fundamental solution as a *convergent* in the continued fraction expansion of \sqrt{d} , and this is relatively fast – it depends upon the *period length*. Alas we shall not touch the theory of continued fractions in this course. If you are interested enough to want to explore the matter further, you will find that the notion of a continued fraction is easy (and fun) to learn (I would rather spend course time discussing things which are not so easy to pick up on one's own).

Continued fractions are not the last word on solving the Pell Equation, however. When d is truly large, other methods are required. Amazingly, a test case for this can be found in the mathematics of antiquity: the so-called **cattle problem of Archimedes**. Archimedes composed a lengthy poem (“twenty-two Greek elegiac distichs”) which is in essence the hardest word problem in human history. The first part, upon careful study, reduces to solving a linear Diophantine equation (in several variables), which is essentially just linear algebra, and it turns out that there is a positive integer solution. However, to get this far is “merely competent”, according to Archimedes. The second part of the problem poses a further constraint which boils down to solving a Pell equation with $d = 410286423278424$. In 1867 the German mathematician C.F. Meyer set out to solve the problem (by hand, of course) using continued fractions. However, he computed 240 steps of the continued fraction expansion of \sqrt{d} , whereas the period length is in fact 203254. Only in 1880 was the problem solved, by A. Amthor. (The gap between the problem and the solution – 2000 years and change – makes the case of Fermat's Last Theorem look fast!) Amthor used a different method. All of this and much more is discussed in

a truly beautiful recent article by Hendrik Lenstra: see

<http://www.ams.org/notices/200202/fea-lenstra.pdf>.

7. SOME FURTHER COMMENTS

There is much more to be said on the subject. Just to further scratch the surface:

It is a purely algebraic consequence of our main result that the unit group of the ring $\mathbb{Z}[\sqrt{d}]$ (for d positive and nonsquare, as usual) is isomorphic to $\mathbb{Z} \times \mathbb{Z}/2\mathbb{Z}$. Indeed, in solving the Pell equation, we found that the group of all norm one units is of this form, and it remains to account for units of norm -1 . Sometimes there are none – e.g. when d is a prime which is $3 \pmod{4}$ – and in this case the result is clear. But in any case the map $N : \mathbb{Z}[\sqrt{d}]^\times \rightarrow \{\pm 1\}$ has as its kernel the solutions to the Pell equation, so if there are also units of norm -1 the units of norm 1 form an index 2 subgroup. On the other hand units of finite order are necessarily roots of unity, of which there are no more than ± 1 in all of \mathbb{R} , let alone $\mathbb{Q}(\sqrt{d})$. The result follows from these considerations; the proof of this is left as an optional exercise.

This is a special case of an extremely important and general result in algebraic number theory. Namely, one can consider any *algebraic number field* – a finite degree field extension K of \mathbb{Q} – and then the ring O_K of all algebraic integers of K – that is, elements α of K which satisfy a monic polynomial with \mathbb{Z} coefficients. We have been looking at the case $K = \mathbb{Q}(\sqrt{d})$, a real quadratic field. Other relatively familiar examples are the cyclotomic fields $\mathbb{Q}(\zeta_N)$ obtained by adjoining an N th root of unity: one can show that this field has degree $\varphi(N)$ over \mathbb{Q} (equivalently, the cyclotomic polynomial Φ_N is irreducible over \mathbb{Q}).

Dirichlet's unit theorem asserts that the units O_K^\times form a finitely generated abelian group, i.e., are isomorphic to $\mathbb{Z}^a \times F$, where F is a finite group (which is in fact the group of roots of unity of F). Noting that the unit group is finite for imaginary quadratic fields and infinite for real quadratic fields, one sees that the rank a must depend upon more than just the degree $d = [K : \mathbb{Q}]$ of the number field: somehow it depends upon “how real” the field is. More precisely, let r be the number of field homomorphisms from K into \mathbb{R} . Alternately, one can show that K is obtained by adjoining a single algebraic number α , i.e., $K = \mathbb{Q}[t]/(P(t))$, where P is a polynomial of degree $d = [K : \mathbb{Q}]$. Then r is nothing else than the number of real roots of the defining (“minimal”) polynomial $P(t)$. In particular $r \leq d$, and $d - r$, the number of complex roots, is even. Then the precise form of Dirichlet's Unit Theorem asserts that $a = r + \frac{d-r}{2} - 1$, a quantity which is positive in every case except for $K = \mathbb{Q}$ and K an imaginary quadratic field! However the proof in the general case requires different techniques.

However, the argument that we used to find the general solution to the Pell equation is fascinating and important. On the face of it, it is very hard to believe that the problem of finding good rational approximations to an irrational number (a problem which is, let's face it, not initially so fascinating) can be used to solve Diophantine equations: we managed to use a result involving real numbers and inequalities to prove a result involving equalities and integers! This is nothing less

than an entirely new tool, lying close to the border between algebraic and analytic number theory (and therefore helping to ensure a steady commerce between them). For instance, in the early 20th century Siegel used Diophantine approximation to prove that a cubic equation of the form $y^2 = x^3 + ax + b$ has only finitely many integral solutions, and in the early 1990's Paul Vojta and Gerd Faltings used it to prove the Mordell-Lang conjecture, which is arguably the single greatest theorem in Diophantine geometry.³ However the subject is also a notoriously difficult one. In fact I know of only one other easy result.

Namely, given an irrational number α , we proved that there are infinitely many rational approximations $\frac{p}{q}$ with $|\alpha - \frac{p}{q}| < \frac{1}{q^2}$. It is natural to ask whether one can do better: can we replace 2 with a higher power of q ?

Certainly we can for some numbers. For instance, define

$$\mathcal{L} = \sum_{n=0}^{\infty} 10^{-k!},$$

the idea here being that we have a decimal expansion in which each lonely 1 is followed by a very long succession of zeros. It is easy to see that the sequence of rational numbers afforded by the partial sums, i.e., $\frac{p_N}{q_N} = \sum_{n=0}^N 10^{-k!}$ give fantastically good approximations: for any positive constants A and B one has

$$|\mathcal{L} - \frac{p_N}{q_N}| < \frac{A}{q_N^B}$$

for all sufficiently large N . On the other hand, Liouville proved the following:

Theorem 8. *Suppose α satisfies a polynomial equation $a_d x^d + \dots + a_1 x + a_0$ with integral coefficients. Then there exists a positive constant A such that for all integers p and $0 \neq q$,*

$$|\alpha - \frac{p}{q}| > \frac{A}{q^d}.$$

That is, being algebraic of degree d imposes an upper limit on the goodness of the approximation by rational numbers. An immediate and striking consequence is that Liouville's number \mathcal{L} cannot satisfy an algebraic equation of any degree: that is, it is a transcendental number! (In fact, by this argument Liouville established the existence of transcendental numbers for the first time.)

Liouville's theorem was improved by many mathematicians, including Thue and Siegel, and culminating in the following theorem of Klaus Roth:

Theorem 9. *(Roth, 1955) Let α be an algebraic real number (of any degree), and let $\epsilon > 0$ be given. Then there are at most finitely many rational numbers $\frac{p}{q}$ satisfying*

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

For this result Roth won the Fields Medal in 1958.

³It therefore pains me not to be able to tell you what it is, but it is not so elementary to state. However it immediately implies Faltings' 1983 theorem that a plane curve of degree at least 4 has only finitely many rational points, which was probably the previous front-runner for the "best theorem ever."