

1. MOTIVATION

Let $d \in \mathbf{Z}$ be a nonsquare positive integer. We have seen that the Pell equation $x^2 - dy^2 = 1$ is closely tied up with the problem of finding units in the ring of integers of the number field $K = \mathbf{Q}(\sqrt{d})$, modulo the minor problem that $\mathbf{Z}[\sqrt{d}]$ may just be an order in \mathcal{O}_K . More specifically, at least when d is squarefree we saw that $\mathbf{Z}[\sqrt{d}]^\times$ is either equal to \mathcal{O}_K^\times or at worst has index 3. A variant on that argument shows that in general (allowing d to have some nontrivial square factors, but still not be a perfect square) $\mathbf{Z}[\sqrt{d}]^\times$ has finite index in \mathcal{O}_K^\times , and in particular it is infinite cyclic up to a sign. Indeed, if $\mathbf{Z}[\sqrt{d}]$ has index f in \mathcal{O}_K then $\mathbf{Z}[\sqrt{d}] = \mathbf{Z} + f\mathcal{O}_K$ (by the general description of quadratic orders in terms of their index in their integral closure) and so $\mathbf{Z}[\sqrt{d}]$ is the full preimage of $\mathbf{Z}/f\mathbf{Z}$ in $\mathcal{O}_K/(f)$. In particular, since $(\mathbf{Z}/f\mathbf{Z})^\times$ is a subgroup of $(\mathcal{O}_K/(f))^\times$, $\mathbf{Z}[\sqrt{d}]^\times$ contains the kernel of the reduction homomorphism $\mathcal{O}_K^\times \rightarrow (\mathcal{O}_K/(f))^\times$, so $\mathbf{Z}[\sqrt{d}]^\times$ has finite index in \mathcal{O}_K^\times that divides the cardinality of $(\mathcal{O}_K/(f))^\times$. (In case d is squarefree we have that $f = 1$ or $f = 2$, and this recovers our earlier analysis of Pell's equation in terms of quadratic units when d is squarefree.)

Writing $\mathbf{Z}[\sqrt{d}]^\times = \langle -1 \rangle \times \varepsilon^{\mathbf{Z}}$, the choice of “fundamental unit” ε of the order $\mathbf{Z}[\sqrt{d}]$ is unique up to sign and inversion. These latter two operations do not affect the value of $N_{K/\mathbf{Q}}(\varepsilon) \in \mathbf{Z}^\times = \{\pm 1\}$, and so if the common norm of the fundamental units of $\mathbf{Z}[\sqrt{d}]$ turns out to be -1 (corresponding to \mathbf{Z} -solutions to $x^2 - dy^2 = -1$) then we pass to powers of ε^2 (up to a sign) to get the solutions to the Pell equation $x^2 - dy^2 = 1$.

But what about the *generalized Pell equation* $x^2 - dy^2 = n$ with nonzero $n \in \mathbf{Z}$? We know that even the case $n = -1$ is rather subtle, being tied up with the issue of whether the fundamental units of \mathcal{O}_K^\times have norm 1 or -1 , a matter for which there is no good general theory at all (i.e., one only knows how to deal with it on a case-by-case basis, by actually finding a fundamental unit and computing its norm). So there is no chance of having a general affirmative existence result as nice as for $n = 1$. Nonetheless, it turns out that by *using* fundamental units, we can give a practical algorithm to determine whether or not $x^2 - dy^2 = n$ has a \mathbf{Z} -solution. This is what is explained and illustrated in this handout. First, we record a general fact which is an easy application of the unit theorem in the real quadratic case:

Theorem 1.1. *Let $d > 0$ be a non-square in \mathbf{Z} . For a nonzero $n \in \mathbf{Z}$, the equation $x^2 - dy^2 = n$ has at most finitely many equivalence classes of solutions in \mathbf{Z} , where we consider two solutions (a, b) and (a', b') to be equivalent if $a + b\sqrt{d} = (a' + b'\sqrt{d})u$ for $u \in \mathbf{Z}[\sqrt{d}]^\times$ with $N_{\mathbf{Q}(\sqrt{d})/\mathbf{Q}}(u) = 1$.*

Proof. Let $K = \mathbf{Q}(\sqrt{d})$. Keep in mind that $\mathbf{Z}[\sqrt{d}]$ is merely an order in \mathcal{O}_K , and that its index in \mathcal{O}_K may be rather large (if d has nontrivial square factors). If (a, b) is a solution to $x^2 - dy^2 = n$ in \mathbf{Z} then the (visibly nonzero) ideal $(a + b\sqrt{d})$ in \mathcal{O}_K has norm $|n|$, and each prime factor of $|n| > 0$ lies beneath at most finitely many (in fact, at most two) prime ideals of \mathcal{O}_K . Hence, there are only finitely many possible prime ideals \mathfrak{p} in \mathcal{O}_K which divide $(a + b\sqrt{d})$ (namely, those \mathfrak{p} over the prime factors of $|n|$), and the multiplicity of each such prime ideal in $(a + b\sqrt{d})$ is bounded in terms of the multiplicity of the corresponding prime factor of $|n|$. We conclude that there are only finitely many possibilities for the ideal $(a + b\sqrt{d})$ in \mathcal{O}_K . If (a, b) and (a', b') are equivalent solutions in the sense defined above then the ideals $(a + b\sqrt{d})$ and $(a' + b'\sqrt{d})$ in \mathcal{O}_K are the same, but conversely if these two principal ideals coincide then we can only conclude that $a + b\sqrt{d} = (a' + b'\sqrt{d})u$ for some $u \in \mathcal{O}_K^\times$.

However, the group of units in $\mathbf{Z}[\sqrt{d}]^\times$ with norm 1 has finite index in \mathcal{O}_K^\times , so the equivalence relation in the theorem is only a mild refinement on the relation that $a + b\sqrt{d}$ and $a' + b'\sqrt{d}$ generate the same ideal in \mathcal{O}_K . In particular, the desired finiteness for the set of equivalence classes of solutions is an immediate consequence of the finiteness for the number of possibilities for the ideal $(a + b\sqrt{d})$ in \mathcal{O}_K and the finiteness of the index in \mathcal{O}_K^\times of the group of norm-1 units in the order $\mathbf{Z}[\sqrt{d}]$. ■

The above was an essentially elementary consequence of the unit theorem for real quadratic fields. Much more serious is figuring out whether there are any nontrivial solutions at all! In the cases $n = \pm 1$ we know

how to find all solutions (and in particular if any exist when $n = -1$), either by the method explained in §4.6 of the text or by more efficient classical methods resting on continued fractions. But when $|n| > 1$ we need a new idea. The key insight is to exploit the logarithm map from the proof of the Dirichlet unit theorem to show that if $u \in \mathbf{Z}[\sqrt{d}]^\times$ is any norm-1 unit distinct from ± 1 (which always exists and which we know how to find) then every equivalence class of solutions to $x^2 - dy^2 = n$ in \mathbf{Z} meets a specific bounded region in \mathbf{R}^2 that is *specified in terms of d , n , and u* , in which case we can then do a brute-force search of integer points in that region (combined with some cleverness to cut down the time of the search) to find representatives of all of the finitely many such equivalence classes.

Observe, by the way, that if (a, b) and (a', b') are two \mathbf{Z} -solutions to $x^2 - dy^2 = n$ then the ratio

$$\frac{a + b\sqrt{d}}{a' + b'\sqrt{d}} \in \mathbf{Q}(\sqrt{d})^\times$$

has norm 1, and so these solutions are equivalent if and only if this ratio lies in $\mathbf{Z}[\sqrt{d}]^\times$. In other words, equivalence amounts to this ratio and its reciprocal both lying in $\mathbf{Z}[\sqrt{d}]$, and that is trivial to check by simply “rationalizing the denominator” and seeing if the resulting new expressions for the ratios in $\mathbf{Q}[\sqrt{d}]$ both lie in $\mathbf{Z}[\sqrt{d}]$. (In fact it suffices to check one of the two ratios; do you see why?)

2. THE BOUNDED REGION

Now we can finally come to the main result which renders effective the task of finding all equivalence classes of solutions to $x^2 - dy^2 = n$ in \mathbf{Z} , including the possibility that there are no solutions at all!

Theorem 2.1. *Let $d \in \mathbf{Z}$ be a non-square positive integer and choose a nonzero $n \in \mathbf{Z}$. Consider the positive square root $\sqrt{d} \in \mathbf{R}$, and choose $u \in \mathbf{Z}[\sqrt{d}]^\times - \{\pm 1\}$ with $N(u) = 1$ and $u > 1$. (We can always arrange $u > 1$ by using negation and/or inversion on an initial choice of non-trivial norm-1 unit if necessary.) Every solution to $x^2 - dy^2 = n$ in \mathbf{Z} is equivalent to one in the bounded region of \mathbf{R}^2 defined by*

$$|x| \leq \frac{\sqrt{|n|}(1 + \sqrt{u})}{2}, \quad |y| \leq \frac{\sqrt{|n|}(1 + \sqrt{u})}{2\sqrt{d}}.$$

Beware that u may be quite complicated: when written in terms of the basis $\{1, \sqrt{d}\}$, its \mathbf{Q} -coefficients may be very large relative to d , even if u is a fundamental unit. But the only way to see if this happens is to actually find an explicit u . Note also that in the statement and proof of the theorem it is not necessary to take u to be a fundamental “norm-1” unit. However, it would be wasteful not to use a fundamental “norm-1” unit, as otherwise we would be searching through a much larger domain than is really necessary.

Proof. Let $K = \mathbf{Q}(\sqrt{d}) \subseteq \mathbf{R}$. We use the logarithmic embedding map $L : K^\times \rightarrow (\mathbf{R}^\times)^2 \subseteq \mathbf{R}^2$ defined by $L(\alpha) = (\log |\alpha|, \log |\bar{\alpha}|)$ where $\alpha \mapsto \bar{\alpha}$ is the nontrivial element of $\text{Gal}(K/\mathbf{Q})$. Obviously $L(\alpha\beta) = L(\alpha) + L(\beta)$. Also, if $v \in \mathcal{O}_K^\times$ then $\log |v| + \log |\bar{v}| = \log |N_{K/\mathbf{Q}}(v)| = \log |\pm 1| = 0$, so $L(\mathcal{O}_K^\times)$ is contained in the line $\{x + y = 0\}$ in \mathbf{R}^2 . In particular, since $u\bar{u} = N_{K/\mathbf{Q}}(u) = 1$, we have $\bar{u} = 1/u$ with $u > 1$ so

$$L(u) = (\log u, \log(1/u)) = \log(u) \cdot (1, -1) \in \mathbf{R}^2$$

with $\log(u) \neq 0$, so $L(u)$ and $(1, 1)$ are a basis of \mathbf{R}^2 . Hence, for any $\alpha \in K^\times$ we can uniquely write

$$L(\alpha) = c_1(\alpha) \cdot (1, 1) + c_2(\alpha) \cdot L(u)$$

for $c_j(\alpha) \in \mathbf{R}$. To determine the first of these two coefficients in terms of α , we bust out coordinates:

$$(c_1(\alpha) + c_2(\alpha) \log(u), c_1(\alpha) - c_2(\alpha) \log(u)) = L(\alpha) = (\log |\alpha|, \log |\bar{\alpha}|),$$

so by adding the coordinates we get

$$2c_1(\alpha) = \log |\alpha\bar{\alpha}| = \log |N_{K/\mathbf{Q}}(\alpha)|.$$

Hence, for a unique $c_2(\alpha) \in \mathbf{R}$ we have

$$L(\alpha) = \frac{1}{2} \cdot \log |N_{K/\mathbf{Q}}(\alpha)| \cdot (1, 1) + c_2(\alpha) \cdot (\log u) \cdot (1, -1)$$

where $N_{K/\mathbf{Q}}(\alpha) \in \mathbf{Q}^\times$.

We now consider $\alpha \in \mathcal{O}_K - \{0\}$, so $N_{K/\mathbf{Q}}(\alpha) \in \mathbf{Z} - \{0\}$, and assume $N_{K/\mathbf{Q}}(\alpha) = n$. Let k be the integer nearest to $c_2(\alpha)$ (whatever it may be), so $c_2(\alpha) = k + \delta$ with $|\delta| \leq 1/2$. Hence,

$$L(\alpha) = \frac{1}{2} \cdot \log |n| \cdot (1, 1) + k \cdot (\log u) \cdot (1, -1) + \delta \cdot (\log u)(1, -1).$$

But $k \cdot (\log u) \cdot (1, -1) = L(u^k)$, so

$$L(\alpha u^{-k}) = \frac{1}{2} \cdot \log |n| \cdot (1, 1) + \delta \cdot (\log u) \cdot (1, -1).$$

Finally, assuming $\alpha = a + b\sqrt{d} \in \mathbf{Z}[\sqrt{d}]$, so this corresponds to a \mathbf{Z} -solution (a, b) to $x^2 - dy^2 = n$, we can run through the preceding considerations and replace α with the equivalent solution αu^{-k} (which preserves the value n of the norm but replaces $L(\alpha)$ with $L(\alpha) - kL(u)$) to stay within the same equivalence class but get to the case where $k = 0$, which is to say $|c_2(\alpha)| \leq 1/2$. Geometrically, scaling by powers of u has the effect of sliding a point $(\alpha, \bar{\alpha})$ in $(\mathbf{R}^\times)^2$ along the hyperbola $x^2 - dy^2 = n$, and we are doing this to get into a domain which is bounded in $(\mathbf{R}^\times)^2$ away from the coordinate axes and away from “ ∞ ”. That is, we claim we are in a compact subspace, and we will make this more explicit by finding bounds after applying the logarithm map $(v, v') \mapsto (\log |v|, \log |v'|)$ which only loses sign information.

With $|c_2(\alpha)| \leq 1/2$, comparing x -coordinates (resp. y -coordinates) in the expressions for $L(\alpha)$ gives

$$\log |a \pm b\sqrt{d}| = \frac{1}{2} \cdot \log |n| \pm c_2(\alpha) \log(u) \leq \frac{1}{2} \cdot (\log |n| + \log u) = \log \sqrt{|n|u}$$

(using that $n = N_{K/\mathbf{Q}}(\alpha) = a^2 - db^2$ and $\log u > 0$), so $|a \pm b\sqrt{d}| \leq \sqrt{|n|u}$. But actually we can do better: since $\log(u) > 0$ and for one of the two signs we have $\pm c_2(\alpha) \leq 0$, for one of the two signs we have $\log |a \pm b\sqrt{d}| \leq \log \sqrt{|n|}$. Thus, as an upper bound on $|a \pm b\sqrt{d}|$ we have $\sqrt{|n|u}$ for one sign and $\sqrt{|n|}$ for the other sign. Using this,

$$|x| = \left| \frac{(a + b\sqrt{d}) + (a - b\sqrt{d})}{2} \right| \leq \frac{\sqrt{|n|}(1 + \sqrt{u})}{2}$$

and

$$|y| = \left| \frac{(a + b\sqrt{d}) - (a - b\sqrt{d})}{2\sqrt{d}} \right| \leq \frac{\sqrt{|n|}(1 + \sqrt{u})}{2\sqrt{d}}$$

■

3. EXAMPLES

We apply the main theorem in two ways: to find all solutions up to equivalence in a case where solutions exist, and to prove there are no \mathbf{Z} -solutions in a case where congruential methods *cannot* be used.

Example 3.1. Consider $x^2 - 15y^2 = 34$. The fundamental unit in $\mathbf{Z}[\sqrt{15}]$ is $u = 4 + \sqrt{15}$, whose norm is 1, so it generates (up to a sign) the group of norm-1 units. By the main theorem, any solution is equivalent to a solution (a, b) with

$$|a| \leq \sqrt{34}(1 + \sqrt{u})/2 \approx 11.1, \quad |b| \leq \sqrt{34}(1 + \sqrt{u})/(2\sqrt{15}) \approx 2.9,$$

so let's consider all possible b 's for the search (as that bound is smaller). For $-2 \leq b \leq 2$, we must check if $34 + 15b^2$ is a perfect square. This happens precisely for $b = \pm 1$. The corresponding value for a up to a sign is 7. In other words, up to a sign every solution to $x^2 - 15y^2 = 34$ is equivalent to one of

$$(7, \pm 1).$$

and these two are inequivalent. For example, the solution $(97, 25)$ turns out to be equivalent to $(7, -1)$ (not to $(7, 1)$!):

$$97 + 25\sqrt{15} = (7 - \sqrt{15})(31 + 8\sqrt{15})$$

with $31 + 8\sqrt{15} \in \mathbf{Z}[\sqrt{15}]$ having norm 1 (as we know it must).

Example 3.2. Consider $x^2 - 82y^2 = 31$. This equation has solutions in \mathbf{Q} , such as $(101/3, 11/3)$ and $(149/11, 15/11)$, and by using either of these it can be shown by a classical geometric method that there are then infinitely many solutions in \mathbf{Q} . (In general, any “smooth” conic over an infinite field with a point over that field must have infinitely many points over that field, and they can even be parameterized in an elegant manner.) But remarkably there is no \mathbf{Z} -solution! The existence of \mathbf{Q} -solutions (and the fact that any rational number lies in $\mathbf{Z}_{(p)}$ for all but finitely many primes p) essentially shows that one cannot hope to use congruential methods to rule out the existence of a \mathbf{Z} -point. An argument with the Chinese Remainder Theorem and the two preceding rational solutions with relatively prime denominators (3 and 11) implies that $x^2 - 82y^2 \equiv 31 \pmod{m}$ has solutions for *every* $m \geq 1$, so indeed the non-existence of \mathbf{Z} -solutions is quite subtle and requires some real technique beyond the kind of elementary congruence arguments which rule out \mathbf{Z} -solutions to $x^2 - 5y^2 = 2$ (look mod 5).

We will use the method of the main theorem, which rests crucially on the existence of non-trivial norm-1 units in real quadratic rings, a fact lying rather deeper than elementary congruential methods. In the ring $\mathbf{Z}[\sqrt{82}]$ a fundamental unit is $9 + \sqrt{82}$, whose norm is -1 , so its square $u = 163 + 18\sqrt{82}$ is a fundamental norm-1 unit. By the main theorem, if $x^2 - 82y^2 = 31$ is to have any \mathbf{Z} -solution at all, it must have such a solution (a, b) with

$$|a| \leq \sqrt{31}(1 + \sqrt{u})/2 \approx 53.05, \quad |b| \leq \sqrt{31}(1 + \sqrt{u})/(2\sqrt{82}) \approx 5.9.$$

To do a search in this box for possible \mathbf{Z} -solutions, we may restrict our attention to the quadrant $a, b \geq 0$. For $0 \leq b \leq 5$ one checks by inspection that $31 + 81b^2$ is never a perfect square! Hence, there are no \mathbf{Z} -solutions, as claimed.