

Figure 4.

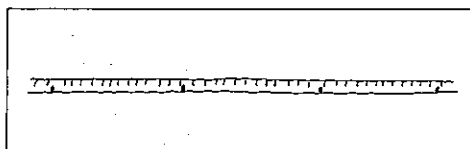


Figure 5. Tally marks showing lunar months and solar years.

other than the derived sequence of the cutting sequence of L relative to Λ . This derived sequence s' , being itself a cutting sequence, is also almost constant. We can compute its value by observing that s' is, of course, the same as the cutting sequence of $\Phi^{-1}(L)$ relative to Λ . Now $\Phi^{-1}(L)$ has slope $\lambda - n_0$, and $\lambda - n_0 < 1$. Thus in s' the roles of a and b are interchanged and b' is isolated. Reflecting in the line $x = y$ interchanges a' and b' and gives a line of slope $1/\lambda - n_0 = [n_1, n_2, \dots]$, from which point the argument repeats.

Does every characteristic sequence occur as the cutting sequence of a line? Not quite; for example, the sequence $b^* a b^*$ has an unfortunate "blip" in the middle. However, it is easy to see that every finite characteristic sequence is linear, that is, it comes from a line, for the sequence of derivations eventually terminates in a single symbol a^n or b^n which is obviously the cutting sequence of a line relative to some derived grid. Applying in succession the inverse derivations we obtain a line segment with the given sequence as its cutting sequence.

Characteristic sequences are nothing other than the limits of linear ones.

Lunar Cycles

Let us pause for a moment and digress to that most ancient of sciences, astronomy. The patterns of occurrences of one heavenly event relative to another, patterns which must surely have been observed from earliest times, provide natural examples of our cutting sequences. For example, in some years twelve new moons would have been observed, in others thirteen. One could well imagine this data recorded by a sequence of tallies along a rod, perhaps as in Figure 5. What more natural question to ask than what is the pattern of tallies which appear? Of course, the anomalies, or irregularities of the heavens, mean that in fact the interval between two like events is never exactly fixed, so that the tally sequence would deviate slightly from any cutting sequence based on two fixed lengths. A calendar based on the assumption of equal intervals would gradually drift away from observation. Nevertheless, David Fowler has speculated that Plato and Eudoxus might have studied the theoretical properties of tally sequences, and perhaps even the problem of relating tally sequences to continued fractions. This is not so unlikely as it sounds when one recalls that the procedure for expressing a number as a continued fraction is closely related to the Euclidean algorithm. The reciprocal subtraction process used in the algorithm was called by the Greeks *anthyphairesis*, and is thought by David Fowler to be the basis of a pre-Eudoxan theory of proportion [5].

Some ancient calendars in fact embody astonishingly accurate astronomical data. For example, in the calendar called the Metonic cycle, found in Babylonia from around 490 B.C. and introduced to Athens by Meton in 432 B.C., one finds the approximation 19 years = 235 months = 6940 days. This gives a mean synodic month of 29.5319 days, compared to the modern value of 29.5305 days. Incidentally, the number 19 is to be found at the back of the Book of Common Prayer in the formula for calculating the date of Easter, and reaches us via the Jewish calculations for Passover. The ratio 19:235 was used in the gearing of the Antikythera Mechanism, a remarkable clockwork calendar dating from about 80 B.C. It can in fact be derived from much cruder data than that in the relevant tally sequence and the continued fraction method.

The Punctured Torus

Leaving the Greeks to their anthyphaireses, let us move on some 2,000 years to hyperbolic geometry. Our original problem has, of course, an analogue in the hyperbolic plane. Taking one of the basic squares

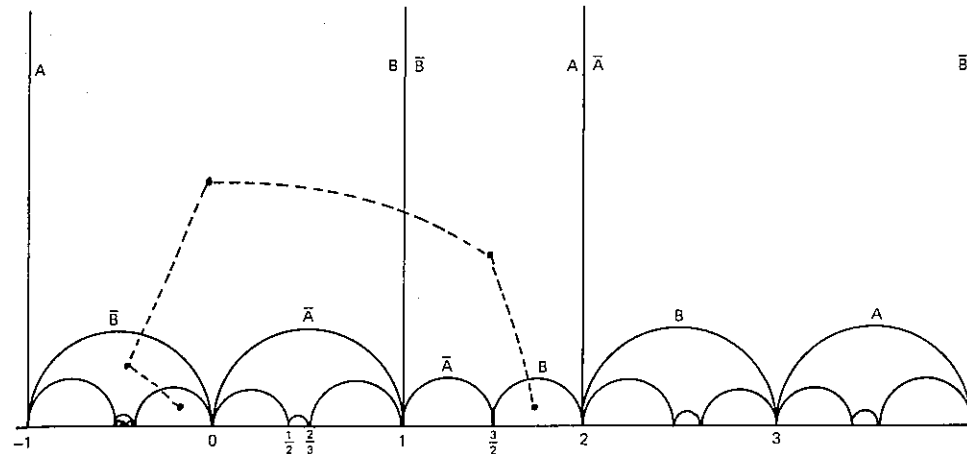


Figure 6. The hyperbolic grid ℓ .

in Λ and glueing the a and b sides, one obtains a torus. We can think of these glueings as implemented by maps $a: (x, y) \rightarrow (x + 1, y)$ and $b: (x, y) \rightarrow (x, y + 1)$. (Incidentally, this explains why we chose to label the sides in Figure 2 as we did.) Now take the hyperbolic grid ℓ illustrated in Figure 6 and glue the A and B sides, this time by the maps $A: z \rightarrow (z + 1)/(z + 2)$ and $B: z \rightarrow (z - 1)/(-z + 2)$.[†] What you get is again a torus, except that since the corners of the "squares" in ℓ are on the boundary of hyperbolic space, one point is missing on the torus and the effect of the hyperbolic metric is to draw out the region around this puncture into an infinitely long spike or cusp as in Figure 1. Just as the maps a, b of the Euclidean plane generate the abelian group \mathbb{Z}^2 which is the fundamental group of the torus, so the maps A, B of the hyperbolic plane generate a free group F which is the fundamental group of the punctured torus T^* . Each "square" in ℓ is an image of the central shaded square S under exactly one element of F , and the labelling of the sides in each square is just a copy of the labels in S .

Recall that straight lines or geodesics in \mathbb{H} are semi-circles centered on \mathbb{R} , or vertical lines. We can pose the same question as before: Which A, B sequences occur as the cutting sequences of lines across ℓ ? Of course, our sequences may now contain not only the symbols A, B but also A^{-1}, B^{-1} (henceforth written as \bar{A}, \bar{B}), depending on the direction in which we cut sides of ℓ .

Observation 3. In a cutting sequence across ℓ a symbol is never immediately followed by its inverse. A se-

quence with this property is called reduced. The solution to our problem is this time remarkably simple: With one exception, every reduced sequence occurs as the cutting sequence of some geodesic in \mathbb{H} , terminating sequences corresponding to lines beginning or ending at the puncture.

The exception is the periodic sequence $\dots \bar{A}\bar{B}\bar{A}\bar{B}\bar{A}\bar{B}\dots$. This corresponds to a loop encircling the puncture, which is a homotopy class with no geodesic representative.

The idea of the proof is to construct a polygonal path in \mathbb{H} whose cutting sequence is the same as that of a given reduced sequence s . This path will consist of line segments joining one square in ℓ to an adjacent one. Each segment is labelled by the side it cuts. Starting from S , we can construct a path whose cutting sequence coincides with s , shown by dotted lines in Figure 6. The fact that s is reduced simply means that the path never doubles back on itself. It is not hard to prove that such paths always converge to two definite distinct points at infinity with the exception of the bad case $(\bar{A}\bar{B}\bar{A}\bar{B})^*$. Joining these points one obtains a geodesic whose cutting sequence is exactly s .

The same method shows that two geodesics have the same cutting sequence if and only if they can be transformed one into the other by an element in F . Since transformations in F preserve ℓ and its labelling, sufficiency is obvious. Suppose that two geodesics have the same cutting sequence. By applying a transformation in F to one of them we can suppose that both cut the same side of ℓ at the same point in their cutting sequence. Fixing an initial side fixes the edge paths of both sequences, which therefore coincide. It follows that the two geodesics are the same.

[†] For more details about hyperbolic geometry and tessellations see the author's earlier *Intelligencer* article "Non-Euclidean Geometry, Continued Fractions and Ergodic Theory" in Vol. 4, No. 1, 1982.

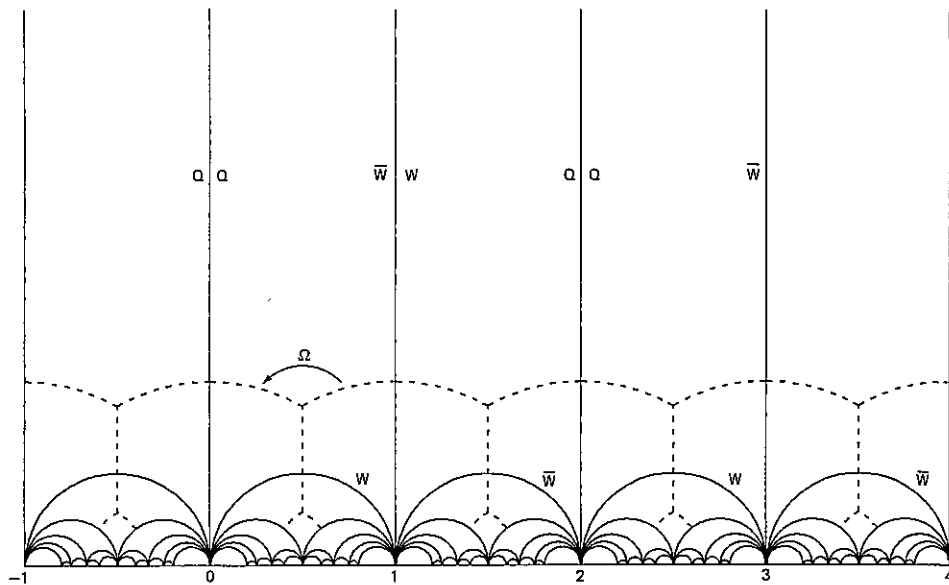


Figure 7. The tessellation Γ subdivided into fundamental regions for $SL(2, \mathbb{Z})$.

$SL(2, \mathbb{Z})$ and Continued Fractions

There was no real reason to take the basic shape S in \mathbb{H} to be a square. One can play the same game with any tessellation \mathcal{T} provided that the vertices of the fundamental region R all lie at infinity. One such tessellation is illustrated in Figure 7. This is associated to $\Gamma(2)$, a subgroup of index 3 in $SL(2, \mathbb{Z})$.[†] The sides of the fundamental region R are mapped to each other by the maps $Q: z \rightarrow -1/z$, $W: z \rightarrow 2 - 1/z$, and this gives the labelling in Figure 7. As shown in the diagram \mathcal{T} is subdivided into three regions each of which is a fundamental region for $SL(2, \mathbb{Z})$. The matrix $\Omega = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$ is a rotation by $2\pi/3$ about $1 + \sqrt{3}i/2$ and rotates these regions onto each other.

The cutting sequences of geodesics relative to \mathcal{T} are of the form $\dots QW^{n_1}QW^{n_2}Q \dots$, where $n_i \in \mathbb{Z}$. Notice that Q^2 never appears since $Q^{-1} = Q$.

Since $SL(2, \mathbb{Z})$ is generated by Q, W and Ω , its action preserves the tessellation \mathcal{T} although not the Q, W labelling. We can, however, label segments of geodesics cutting across the triangles in \mathcal{T} so as to be invariant under $SL(2, \mathbb{Z})$, by labelling a segment L or R according to whether the vertex of the triangle cut off by the

segment is to left or right, as we have done in Figure 8. It is easy to write down a recipe for conversion from Q, W to L, R sequences:

$$\begin{array}{cccccc} QW & \bar{W}\bar{W} & WQ & Q\bar{W} & WW & \bar{W}Q \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ L & L & L & R & R & R \end{array}$$

It now follows that two geodesics in \mathbb{H} are equivalent under $SL(2, \mathbb{Z})$ if and only if their L, R sequences agree.

But this is not all. The L, R sequences bring us back to continued fractions! Let θ be any positive real number, and as in Figure 8 let $\gamma(\theta)$ be a geodesic ray joining any point on the imaginary axis to θ . Reading off the L, R sequence of $\gamma(\theta)$ we obtain a sequence $L^{n_0}R^{n_1}L^{n_2} \dots$ (if $\theta < 1$ the sequence begins with R not L). Then $[n_0, n_1, n_2, \dots]$ is the continued fraction expansion of θ !

The proof is not hard. First, it is obvious that $n_0 = \lfloor \theta \rfloor$. Let D be the point where $\gamma(\theta)$ cuts $\theta = n_0$. Applying the map $\tau_1: z \rightarrow -1/z - n_0$, D is mapped to a point D' on the imaginary axis and $\gamma(\theta)$ becomes a ray γ' through D' pointing in the negative direction with endpoint at $-1/\theta - n_0$. The n_1 segments of type R in $\gamma(\theta)$ which follow the initial n_0 segments of type L are now apparent as the n_1 vertical strips crossed by $\tau_1(\gamma)$ before it descends to $\tau_1(\theta)$. Thus $n_1 = 1/\theta - n_0$, so that $\theta = n_0 + 1/n_1 + r$, $0 < r < 1$. Now apply $\tau_1: z \rightarrow -1/(z + n_1)$ to γ' and proceed as before [9].

[†] Recall $SL(2, \mathbb{Z}) = \{ \begin{pmatrix} a & b \\ c & d \end{pmatrix} \mid a, b, c, d \in \mathbb{Z}, ad - bc = 1 \}$. Of course $SL(2, \mathbb{Z})$ acts on \mathbb{H} by $z \rightarrow az + b/cz + d$.

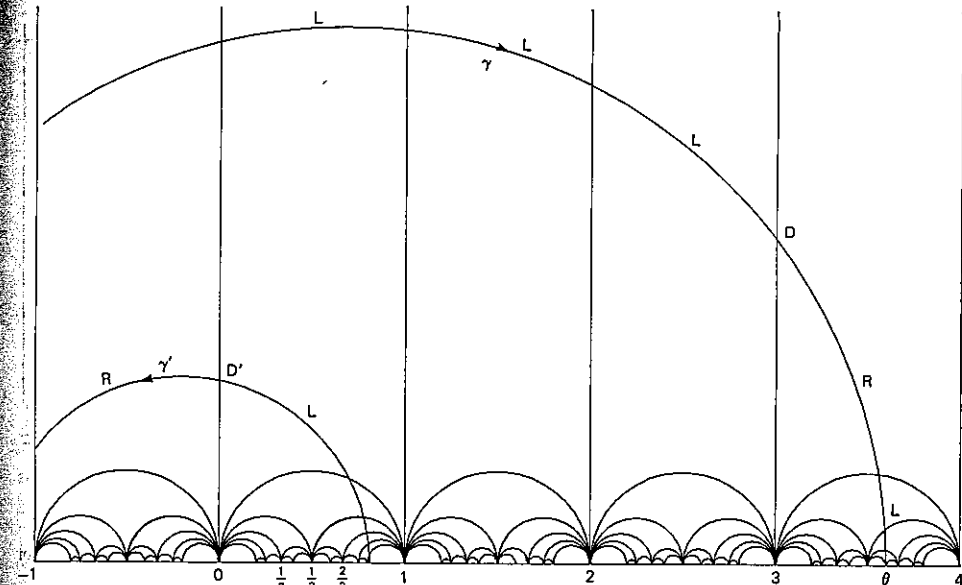


Figure 8. Reading off the continued fraction expansion of θ from \mathcal{T} : $\theta = 3 + \frac{1}{1 + \dots}$.

Simple Curves on the Punctured Torus and the Dickson Rules

We indicated at the beginning that Markoff irrationalities are associated to simple loops on the punctured torus T^* . We are now in a position to understand exactly what these loops are. In fact: A geodesic on T^* is closed and simple if and only if its cutting sequence is periodic and characteristic. By the cutting sequence of a curve on T^* we mean, of course, the cutting sequence of any of its lifts to \mathbb{H} . Since all these lifts are equivalent under F , we know that cutting sequences coming from different lifts are the same. Closed geodesics correspond exactly to those with periodic cutting sequences. We know that periodic characteristic sequences correspond exactly to lines of rational slope on the square grid Λ . Let L be such a line, and move L if necessary so as to avoid the vertices of Λ .

Since L is disjoint from all its images under the vertical and horizontal translations a and b , its image on T^* cannot contain any self-intersections; in other words, it is simple. Now there is exactly one F -equivalence class of geodesics on the hyperbolic plane \mathbb{H} with the same cutting sequence as L , and it is not hard to show that the corresponding geodesic on T^* is also simple. This geodesic is obtained, if you like, by pulling tight the curve L on T relative to the hyperbolic metric on T^* .

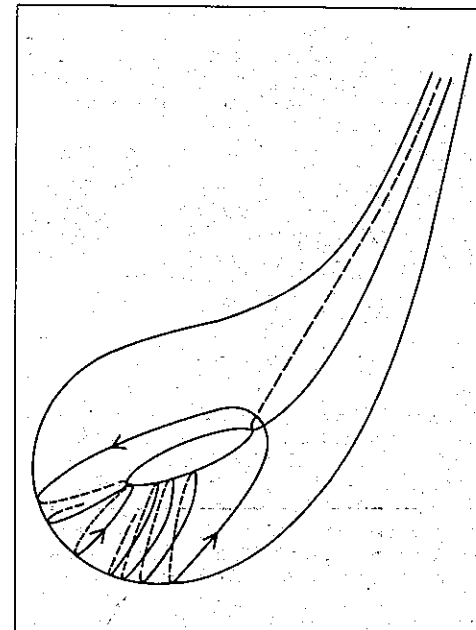


Figure 9. The curve $\dots AAAABAAA \dots$.

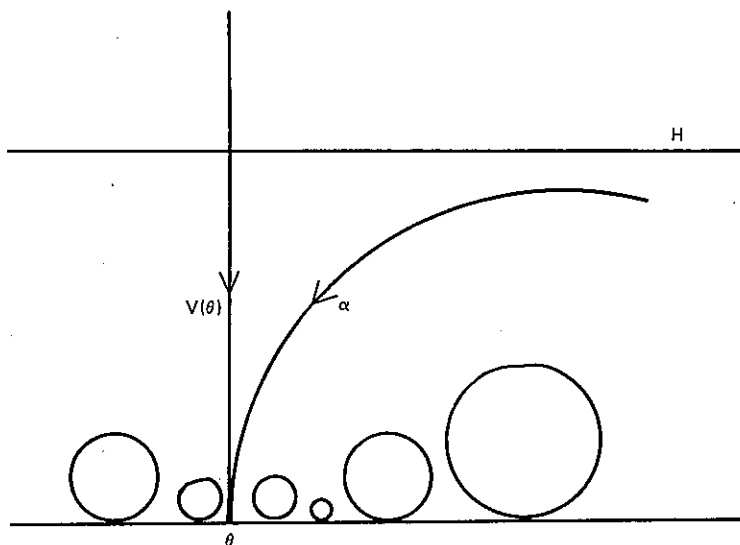


Figure 12. Simple curves avoiding horocycles.

tance outside \bar{H} , we see that $v(\theta)$ enters N only a finite number of times. Thus one sees from Figure 12 that the inequality $|\theta - a/c| < 1/3c^2$ has only a finite number of solutions so that $v(\theta) > 1/3$.

One can actually calculate that the closest approach of α to \bar{H} is $\log \coth l(\alpha)/2$ where $l(\alpha)$ is the hyperbolic length of α [6]. This gives an exact value for $v(\theta)$.

In case (ii) the tail of $v(\theta)$ is characteristic and hence \bar{v} can be approximated by a sequence of simple curves $\bar{\alpha}_n$. Since there are only finitely many curves with lengths below a given bound, the sequence of lengths tends to infinity and hence the distance to \bar{H} approaches zero. Thus $v(\theta)$ approaches N arbitrarily closely although from some point on it never enters N since $s_n(\theta)$ is eventually characteristic. Combining these facts one sees that $v(\theta) = 1/3$.

Finally, in case (iii) the tails $s_n(\theta)$ are never characteristic and so $v(\theta)$ enters N infinitely often. Thus there are infinitely many solutions to $|\theta - a/c| < 1/3c^2$; in other words, $v(\theta) \leq 1/3$.

Trace Formulae, Diophantine Equations and Quadratic Forms

Hoping that the reader's patience is not completely exhausted, we will conclude by giving some brief pointers to the connection of our approach to another well known aspect of Markoff's theory, the minima of binary quadratic forms.

The Markoff spectrum is frequently calculated by introducing Markoff triples. These are integer triples

(x, y, z) which are solutions of the Diophantine equation

$$x^2 + y^2 + z^2 = 3xyz. \quad (D)$$

Associated to such a triple is a pair of real quadratic numbers $\xi, \xi' = 1/2 + y/xz + 1/2(9 - 4z^2)^{1/2}$. The numbers ξ, ξ' are Markoff irrationalities with $v(\xi) = v(\xi') = \sqrt{9 - 4z^2} > 1/3$.

In fact, as explained in [2], Markoff triples are (up to a factor of 3) the traces of triples $(U, V, \bar{V}U)$ such that U, V are a pair of generators for the group F with fundamental region as shown in Figure 13. The simplest solution $(1, 1, 1)$ corresponds to the A, B generators we used above. The formula (D) is nothing other than one of Fricke's trace identities relating traces of matrices in $SL(2, \mathbb{R})$. Starting with the solution $(1, 1, 1)$ the operations $(x, y, z) \rightarrow (z, x, y)$ and $(x, y, z) \rightarrow (x, 3xy - z, y)$ generate all possible solutions to (D). These operations are the same as the operations of derivation and substitution which we used above.

The geodesics $\bar{\gamma}(A), \bar{\gamma}(B)$ corresponding to the minimal solution $(1, 1, 1)$ of (D) are simple. Since the operation of derivation is induced by an isometry of T^* , the same is actually true of all solutions of (D). Now the geodesic $\bar{\gamma}(M)$ associated to a matrix $M = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \in F$ is the projection of a geodesic $\gamma(M)$ on \mathbb{H} whose endpoints are the fixed points ξ_M, ξ'_M of M on \mathbb{R} . These endpoints are of course roots of $c\xi^2 + (d - a)\xi - b = 0$. Thus we see in another way that Markoff irrationalities are the endpoints of lifts of simple geodesics on T^* .

One can associate to M the quadratic form

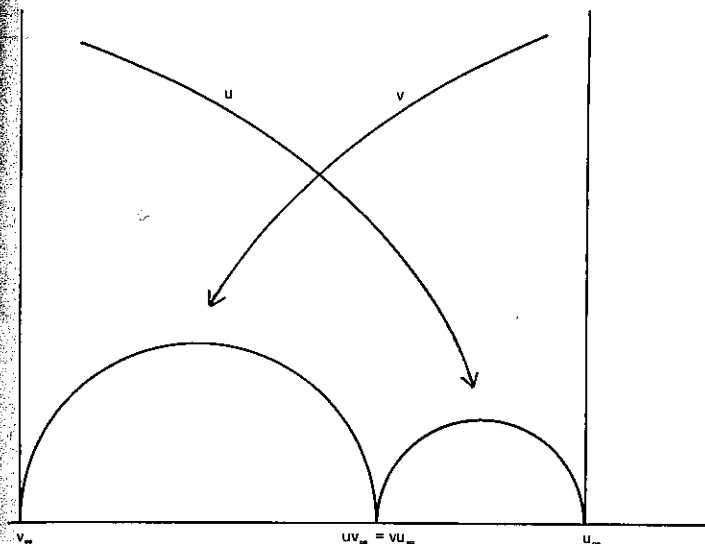


Figure 13. Fundamental region for the punctured torus.

$$Q_M(x, y) = cx^2 + (d - a)xy - by^2.$$

Since $\gamma(M)$ is simple it lies below the line $\text{Im}z = 3/2$, in other words, $|\xi_M - \xi'_M| < 3$. Now $|\xi_M - \xi'_M| = \Delta^{1/2}/Q_M(1, 0)$ where $\Delta = \text{Tr}^2 M - 4$ is the discriminant of Q_M , so that $Q_M(1, 0)/\Delta^{1/2} > 1/3$.

But we know more. Since the action of $SL(2, \mathbb{Z})$ on \mathbb{H} preserves L, R sequences, it preserves simple geodesics on T^* . Thus if $\gamma(M)$ is simple, so is $\gamma(gMg^{-1})$ for any $g \in SL(2, \mathbb{Z})$. One easily computes that

$$Q_M(x, y) = Q_{gMg^{-1}}(gx, gy)$$

and that Q_M and $Q_{gMg^{-1}}$ have the same discriminant Δ . Given any pair $(x, y) \in \mathbb{Z}^2$ we can always find $g \in SL(2, \mathbb{Z})$ with $(gx, gy) = (1, 0)$. Hence

$$Q_M(x, y) = Q_{gMg^{-1}}(1, 0) > \Delta^{1/2}/3;$$

in other words,

$$\min_{x, y \in \mathbb{Z}^2} \frac{Q_M(x, y)}{\Delta^{1/2}} > 1/3.$$

Of course the actual value of the minimum can be calculated and is, not surprisingly, $v(\xi_M)$. For matrices M which do not correspond to simple geodesics, the minimum lies on or below the value $1/3$.

These are the results of Markoff on minima of binary quadratic forms.

It seems clear from the foregoing that the next level of approximation should be studied by looking at geodesics with one self-intersection. Such geodesics penetrate only a bounded distance into \bar{H} . One wonders

if these further levels of approximation are perhaps related to phenomena of successive transitions from periodicity into chaos?

References

1. E. B. Christoffel, *Observatio Arithmetica*, *Annali di Matematica*, 2nd series, 6(1875), 148-152.
2. H. Cohn, Approach to Markoff's minimal forms through modular functions. *Ann. Math.* 61(1955), 1-12.
3. H. Cohn, Representation of Markoff's binary quadratic forms by geodesics on a perforated torus. *Acta Arithmetica* XVIII(1971), 125-136.
4. L. E. Dickson, *Studies in the theory of numbers*. Chicago: 1930.
5. D. Fowler, Anthyphairctic ratio and Eudoxan proportion. *Archive for History of Exact Sciences* 24(1981), 69-72.
6. A. Haas, Diophantine approximation on hyperbolic Riemann surfaces, *Bull. A.M.S.* 11(1984), 359-362.
7. J. Lehner, M. Scheingorn, Simple closed geodesics on $H^*/\Gamma(3)$ arise from the Markoff spectrum, preprint.
8. A. A. Markoff, Sur les formes binaires indefinies, I, *Math. Ann.* 15(1879), 281-309; II, 17(1880), 379-400.
9. C. Series, The modular surface and continued fractions. *J. London Math. Soc.* (1984).
10. A. L. Schmidt, Minimum of quadratic forms with respect to Fuchsian groups I. *J. Reine Angew. Math.* 286/7 (1976), 341-368.
11. H. J. S. Smith, Note on continued fractions. *Messenger of Mathematics*, 2nd series, 6(1876), 1-14.
12. E. C. Zeeman, An algorithm for Eudoxan and anthyphairctic ratios, preprint.

Department of Mathematics
University of Pennsylvania
Philadelphia, Pennsylvania 19104 U.S.A.

