

Ratner's Theorem on Horocyclic Flows

John H. Hubbard and Robyn L. Miller
Cornell University

December 8, 2007

1 Horocycle flow on hyperbolic surfaces

Let X be a complete hyperbolic surface, perhaps the hyperbolic plane H , and let \mathbf{X} denote the unit tangent bundle $T^1(X)$ to X (and $\mathbf{H} = T^1H$). There are three flows on \mathbf{X} which will concern us here. They are realized by three cars, as represented in Figure 1.

The cars all have their steering wheels locked in position: the first car drives straight ahead, the second one steers to the left so as to follow a path of geodesic curvature 1, and the third steers to the right, also following a path of geodesic curvature 1. All three cars have an arrow painted on the roof, centered at the rear axle; for the first the arrow points straight ahead, and for the other two it points sideways – in the direction towards which the car is steering for the second car and in the opposite direction for the third.

The flows at time $t \in \mathbb{R}$ starting at a point $\mathbf{x} = (x, \xi) \in \mathbf{X}$ are defined as follows:

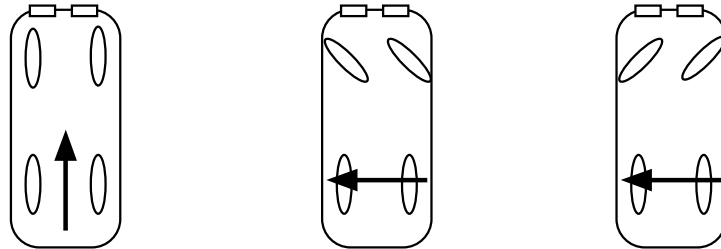


Figure 1: Driving the cars above leads to geodesic flow, positive horocycle flow and negative horocycle flow respectively

1. *The geodesic flow*: put the first car on X with the arrow pointing in the direction of ξ , and drive a distance t . The point of arrival, with the arrow on the car at that point, will be denoted $\mathbf{x}g(t)$;
2. *The positive horocyclic flow*: put the second car on X with the arrow pointing in the direction of ξ , and drive a distance t . The point of arrival, with the arrow on the car at that point, will be denoted $\mathbf{x}u_+(t)$;
3. *The negative horocyclic flow*: put the third car on X with the arrow pointing in the direction of ξ , and drive a distance t . The point of arrival, with the arrow on the car at that point, will be denoted $\mathbf{x}u_-(t)$.

We will see when we translate to matrices why it is convenient to write the flows as *right* actions.

The trajectories followed by these cars are represented in Figure 2.

2 Translation to Matrices

In less picturesque language (more formal, not more accurate), you can identify \mathbf{X} with $\Gamma \backslash \mathrm{SL}_2(\mathbb{R})$ for some Fuchsian group Γ .

1. The geodesic flow of the point represented by $g \in \mathrm{SL}_2(\mathbb{R})$ is

$$t \mapsto g \begin{pmatrix} e^{\frac{t}{2}} & 0 \\ 0 & e^{-\frac{t}{2}} \end{pmatrix};$$

2. The positive horocyclic flow of the point represented by $g \in \mathrm{SL}_2(\mathbb{R})$ is

$$t \mapsto g \begin{pmatrix} 1 & t \\ 0 & 1 \end{pmatrix};$$

3. The negative horocyclic flow of the point represented by $g \in \mathrm{SL}_2(\mathbb{R})$ is

$$t \mapsto g \begin{pmatrix} 1 & 0 \\ t & 1 \end{pmatrix};$$

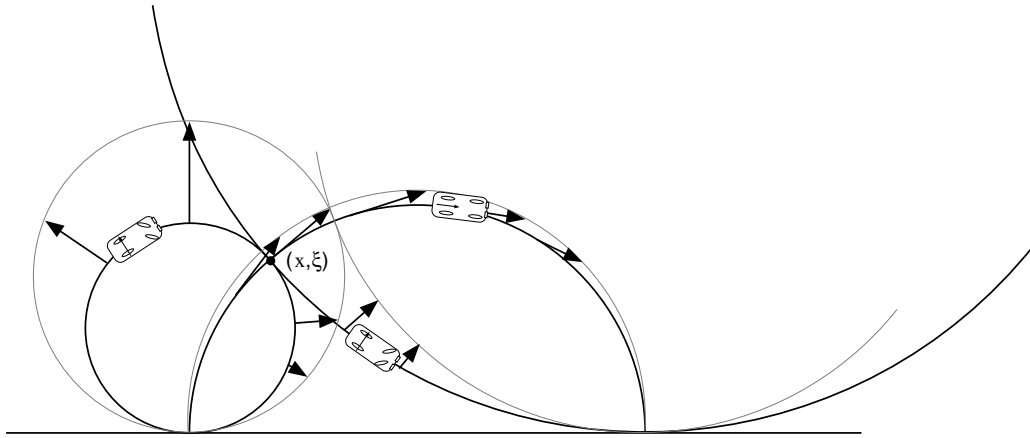


Figure 2: In the upper half-plane model of the hyperbolic plane, the geodesic passing through (x, ξ) is the semicircle perpendicular to the real axis and tangent at x to ξ . One should remember that it is not a curve in \mathbb{H} , but rather a curve in $T^1(\mathbb{H})$ and carries its velocity vector with it. From the point (x, ξ) , the positive horocycle flow is the circle tangent to the real axis at the endpoint of the geodesic above and perpendicular to ξ at x , whereas the negative horocycle flow is the circle tangent to the real axis at the origin of the geodesic, and still perpendicular to ξ at x . We have drawn our tinkertoys driving along them.

The standard left action of $\mathrm{SL}_2(\mathbb{R})$ on H , which lifts by the derivative to a left action on \mathbf{H} is given by

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \cdot z = \frac{az + b}{cz + d} \quad \text{lifting to} \quad \begin{bmatrix} a & b \\ c & d \end{bmatrix} \cdot (z, \xi) = \left(\frac{az + b}{cz + d}, \frac{\xi}{(cz + d)^2} \right) \quad (1)$$

We can then identify $\mathrm{SL}_2 \mathbb{R}$ to \mathbf{H} by choosing $\mathbf{x}_0 = (i, i) \in \mathbf{H}$ and setting $\Phi : \mathrm{SL}_2(\mathbb{R}) \rightarrow \mathbf{H}$ to be

$$\Phi \left(\begin{bmatrix} a & b \\ c & d \end{bmatrix} \right) := \begin{bmatrix} a & b \\ c & d \end{bmatrix} \cdot \mathbf{x}_0 = \left(\frac{ai + b}{ci + d}, \frac{i}{(ci + d)^2} \right)$$

Since

$$\Phi(\gamma A) = (\gamma A) \cdot \mathbf{x}_0 = \gamma \cdot (A \cdot \mathbf{x}_0) = \gamma \cdot \Phi(A)$$

we see that Φ induces a diffeomorphism $\Phi_\Gamma : \Gamma \backslash \mathrm{SL}_2(\mathbb{R}) \rightarrow \Gamma \backslash \mathbf{X}$.

The left action above does *not* induce an action of $\mathrm{SL}_2 \mathbb{R}$ on $\Gamma \backslash \mathbf{X}$, but there is an action on the right given by

$$\Phi(A) * B = \Phi(AB)$$

For $\gamma \in \Gamma$ we have

$$\Phi_\Gamma(A) * B = \Phi_\Gamma(AB) = \Phi_\Gamma(\gamma AB) = \gamma \cdot \Phi_\Gamma(AB) = \gamma \cdot (\Phi_\Gamma(A) * B)$$

so the action is well defined on \mathbf{X} . All three flows are special cases, eg. write

$$G^t = \begin{pmatrix} e^{\frac{t}{2}} & 0 \\ 0 & e^{-\frac{t}{2}} \end{pmatrix}, \quad U_+^t = \begin{pmatrix} 1 & t \\ 0 & 1 \end{pmatrix}, \quad U_-^t = \begin{pmatrix} 1 & 0 \\ t & 1 \end{pmatrix}$$

and name the corresponding one-parameter subgroups

$$G = \{G^t \mid t \in \mathbb{R}\}, \quad U_+ = \{U_+^t \mid t \in \mathbb{R}\}, \quad U_- = \{U_-^t \mid t \in \mathbb{R}\}.$$

Then

$$\mathbf{x}g(t) = \mathbf{x} * G^t, \quad \mathbf{x}u_+(t) = \mathbf{x} * U_+^t, \quad \mathbf{x}u_-(t) = \mathbf{x} * U_-^t$$

Let us check these. By naturality we see that for all $A \in \mathrm{SL}_2(\mathbb{R})$ we have

$$(A \cdot \mathbf{x}_0)g(t) = A \cdot (\mathbf{x}_0g(t)), \quad (A \cdot \mathbf{x}_0)u_+(t) = A \cdot (\mathbf{x}_0u_+(t)), \quad (A \cdot \mathbf{x}_0)u_-(t) = A \cdot (\mathbf{x}_0u_-(t))$$

and, moreover

$$\Phi(G^t) = \mathbf{x}_0g(t), \quad \Phi(U_+^t) = \mathbf{x}_0u_+(t), \quad \Phi(U_-^t) = \mathbf{x}_0u_-(t)$$

so

$$\Phi(AG^t) = (AG^t) \cdot \mathbf{x}_0 = A \cdot (G^t \cdot \mathbf{x}_0) = A \cdot (\mathbf{x}_0 g(t)) = (A \cdot \mathbf{x}_0) g(t) = \Phi(A)g(t)$$

and the argument for u_+ and u_- is identical.

Left multiplication by G^t , U_+^t and $U_-(t)$ also give flows on \mathbf{H} ; probably easier to understand than the geodesic and horocycle flows. For instance, left action by U_+^t corresponds to translating a point and vector to the right by t . But these actions do not commute with the action of Γ and hence induce nothing on \mathbf{X} .

Since $\mathrm{SL}_2 \mathbb{R}$ is unimodular, it has a Haar measure, invariant under both left and right translation, and unique up to multiples. Since $\mathrm{SL}_2 \mathbb{R}$ is not compact, there is no natural normalization. Denote by ω the corresponding measure on \mathbf{H} ; if $\mathbf{X} = \Gamma \backslash \mathbf{H}$ is of finite volume, we will denote by $\omega_{\mathbf{X}}$ the corresponding measure normalized so that $\omega_{\mathbf{X}}(\mathbf{X}) = 1$. Up to a constant multiple we have

$$\omega = \frac{dx \wedge dy \wedge d\theta}{y^2},$$

where we have written $\mathbf{x} = (z, \xi)$ and $z = x + iy$, $\xi = ye^{i\theta}$ (the factor y is there to make it a unit vector): this measure is easily confirmed to be invariant under both left and right action of $\mathrm{SL}_2 \mathbb{R}$ on \mathbf{H} .

Occasionally, we will need a metric and not just a measure on $\mathrm{SL}_2 \mathbb{R}$; we will use the metric that corresponds under Φ to the Riemannian structure

$$\frac{dx^2 + dy^2}{y^2} + d\theta^2.$$

This metric is invariant under left action of $\mathrm{SL}_2 \mathbb{R}$ on \mathbf{H} , and as such does induce a metric on \mathbf{X} . It is *not invariant* under right action, and the flows u_+ , u_- and g do not preserve lengths.

3 The Horocycle Flow is Ergodic

Theorem 1 (Hedlund) [*Hed36*] *The positive and the negative horocycle flows are ergodic.*

We will show this for the positive ergodic flow. To prove Theorem 1 we will show that any $f \in C_c(\mathbf{X})$ invariant under the horocycle flow is constant

almost everywhere. Indeed, if the positive horocycle flow is not ergodic then there is a set $\mathbf{Y} \in \mathbf{X}$ of positive but not full measure that is invariant under U_+ and the characteristic function $\mathbf{1}_{\mathbf{Y}}$ provides a nonconstant invariant function.

Lemma 2 For $g \in L^2(\mathbf{X})$, $A \in \mathrm{SL}_2(\mathbb{R})$ and $\mathbf{x} \in \mathbf{X}$, let $(T_A g)(\mathbf{x}) := g(\mathbf{x} * A)$. Then the function $F_g : \mathrm{SL}_2(\mathbb{R}) \rightarrow \mathbb{R}$ defined by

$$F_g(A) = \int_{\mathbf{X}} g(\mathbf{x})g(\mathbf{x} * A)\omega_{\mathbf{X}}(d\mathbf{x}) := \langle g, T_A g \rangle$$

is

- (a) uniformly continuous and
- (b) bi-invariant under U_+ , i.e., invariant under the left and the right action of U_+ on $\mathrm{SL}_2 \mathbb{R}$.

Proof of Lemma 2 (a) Choose $\varepsilon > 0$. Since the continuous functions with compact support are dense in $L^2(\mathbf{X})$, we can find a function $f \in C_c(\mathbf{X})$ with $\|f - g\|_2 < \varepsilon/3$. The fact that such a f is uniformly continuous means that $\exists \delta > 0$ such that

$$d(A, B) < \delta \Rightarrow \|T_A f - T_B f\|_2 < \frac{\varepsilon}{3}$$

So when $d(A, B) \leq \delta$ we have

$$\|T_A g - T_B g\|_2 \leq \|T_A g - T_A f\|_2 + \|T_A f - T_B f\|_2 + \|T_B f - T_B g\|_2 \leq \varepsilon$$

We see that $A \mapsto T_A g$ is a uniformly continuous map $\mathrm{SL}_2(\mathbb{R}) \rightarrow L^2(\mathbf{X})$ and (a) follows.

(b) Biinvariance reflects the invariance of Haar measure on $\mathrm{SL}_2 \mathbb{R}$ under left and right translation: for $A \in \mathrm{SL}_2(\mathbb{R})$ we have

$$\begin{aligned} F_f(AU_+^t) &= \int_{\mathbf{X}} f(\mathbf{x})f(\mathbf{x} * (AU_+^t))\omega_{\mathbf{X}}(d\mathbf{x}) = \int_{\mathbf{X}} f(\mathbf{x})f((\mathbf{x} * A)u_+(t))\omega_{\mathbf{X}}(d\mathbf{x}) \\ &= \int_{\mathbf{X}} f(\mathbf{x})f(\mathbf{x} * A)\omega_{\mathbf{X}}(d\mathbf{x}) = F_f(A); \\ F_f(U_+^t A) &= \int_{\mathbf{X}} f(\mathbf{x})f(\mathbf{x} * (U_+^t A))\omega_{\mathbf{X}}(d\mathbf{x}) = \int_{\mathbf{X}} f(\mathbf{x} * A^{-1})f(\mathbf{x} * U_+^t)\omega_{\mathbf{X}}(d\mathbf{x}) \\ &= \int_{\mathbf{X}} f(\mathbf{x} * A^{-1})f(\mathbf{x}u_+(t))\omega_{\mathbf{X}}(d\mathbf{x}) = \int_{\mathbf{X}} f(\mathbf{x} * A^{-1})f(\mathbf{x})\omega_{\mathbf{X}}(d\mathbf{x}) \\ &= \int_{\mathbf{X}} f(\mathbf{x})f(\mathbf{x} * A)\omega_{\mathbf{X}}(d\mathbf{x}) = F_f(A) \end{aligned}$$

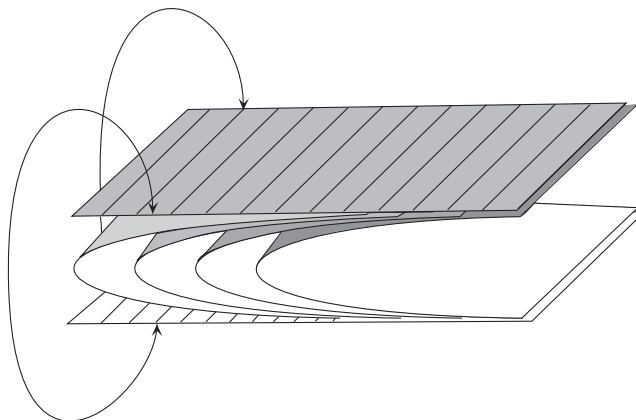


Figure 3: Since the left and the right action of U_+ commute, the bi-orbits are homeomorphic to \mathbb{R}^2 , except the orbits on which the two actions coincide. Viewed in \mathbf{H} , the orbits are the 1-parameter family of folded planes (the ham slices in the sandwich). The top and the bottom planes should be identified; they represent the orbits formed of vertical upwards pointing vectors; those orbits are lines, as drawn in the planes. The salient features of the figure is that any two points of the top plane are within ϵ of a single bi-orbit (in fact all of those with folds sufficiently far to the left), and every bi-orbit comes arbitrarily close to a bi-orbit consisting of vertical upwards pointing vectors.

Proof of Theorem 1. What do the bi-orbits of U^+ look like in $\mathrm{SL}_2 \mathbb{R}$? Using our identification $\Phi : \mathrm{SL}_2 \mathbb{R} \rightarrow \mathbf{H}$, we can think of the bi-orbits as living in \mathbf{H} with geometry represented in figure 3. More specifically, there are two kinds of bi-orbits. The first (exceptional) kind of bi-orbit consists of all upwards pointing vertical vectors anchored at points $z = x + iy$ with a given y -coordinate. The union of these orbits forms a plane $\mathbf{V} \subset \mathbf{H}$.

The other (generic) bi-orbits consist of the vectors defining horocycles of a given radius tangent to the x -axis: each such bi-orbit is diffeomorphic to a plane.

In particular, all the bi-orbits are closed, and there is nothing to prevent the existence of nonconstant continuous functions on \mathbf{H} that are constant on bi-orbits. But our function F is *uniformly* continuous, and that changes the situation: every uniformly continuous function on \mathbf{H} that is constant on bi-orbits *is* constant. What we need to see is that for every $\epsilon > 0$,

- any two elements of \mathbf{V} can be approximated to within ϵ by a single 2-dimensional biorbit, and
- that any 2-dimensional biorbit is within ϵ of some element of \mathbf{V} .

These features are illustrated, but not proved, by figure 3. The proofs are the content of the two parts of figure 4:



Figure 4: Left: Any two upward-pointing vertical unit vectors can be approximated by elements of the same biorbit. Right: Any orbit contains vectors arbitrarily close to upward-pointing vertical.

Suppose that $(z, \xi), (z', \xi') \in \mathbf{V}$ are anchored at $z = x + iy$ and $z' = x' + iy'$ respectively, and that $y > y'$. Take a horocycle passing through z , corresponding to a circle of large radius tangent to the x -axis. Then the defining vector for this horocycle at z is almost vertical, and the defining vector at z' is even more vertical. This is illustrated on the left of figure 4.

The right side of figure 4 shows that the anchored vectors at points sufficiently close to the x -axis on any non-horizontal horocycle are almost in \mathbf{V} .

■

4 Ratner's Theorem

Theorem 3 [Rat92] *Let X be a complete hyperbolic surface of finite area. Then every horocycle on \mathbf{X} is either periodic or equidistributed in \mathbf{X} .*

This theorem is evidently a much stronger statement than that the horocycle flow is ergodic, or even that it is uniquely ergodic. It is not an “almost everywhere” statement, but rather it asserts that *every horocycle* is either periodic or equidistributed in \mathbf{X} .

Note that this depends crucially on the fact that horocycles have geodesic curvature 1. The statement is false for geodesics (geodesic curvature 0): geodesics can do all sorts of things other than being periodic or equidistributed. For instance, they can spiral towards a closed geodesic, or be dense in a geodesic lamination, or spiral towards a geodesic lamination. Curves with constant geodesic curvature < 1 stay a bounded distance away from a geodesic, and hence can do more or less the same things as geodesics; in particular, they do not have to be periodic or equidistributed.

On the other hand, curves with geodesic curvature > 1 are always periodic, hence never equidistributed.

5 The Geometry of Flows in \mathbf{H}

The geodesic flow in \mathbf{X} has *stable* and *unstable* foliations: two points $\mathbf{x}_1, \mathbf{x}_2 \in \mathbf{X}$ belong to the same leaf of the stable foliation if $d(\mathbf{x}_1g(t), \mathbf{x}_2g(t))$ is bounded as $t \rightarrow +\infty$, and they belong to the same leaf of the unstable foliation if $d(\mathbf{x}_1g(t), \mathbf{x}_2g(t))$ is bounded as $t \rightarrow -\infty$. These foliations are very easy to visualize in $T^1\mathbf{H}$, as shown in Figure 5.

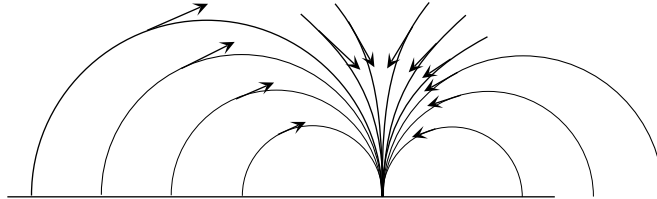


Figure 5: The geodesics all ending at the same point at infinity, together with their tangent vectors, form one leaf of the stable foliation for the geodesic flow. Similarly, the geodesics emanating from a point at infinity form a leaf of the unstable foliation. In \mathbf{X} , these leaves are tangled up in some very complicated way (after all, most geodesics are dense, never mind their stable and unstable manifolds). But clearly each leaf is an immersed smooth surface, hence of measure 0 in \mathbf{X} .

Note that the stable leaves are fixed by the positive horocycle flow: the positive horocycles are the curves orthogonal to the geodesics in a leaf. Similarly, the unstable manifolds are fixed by the negative geodesic flow, and the negative horocycles in a leaf are the curves orthogonal to the geodesics in that leaf.

On the other hand, the positive horocycles are transverse to the unstable manifolds, and positive horocycle flow does not send unstable leaves to unstable leaves.

Set $S_{\mathbf{x}}$ to be the unstable manifold of the geodesic through \mathbf{x} . Define $S_{\mathbf{x}}(a, b) \subset S_{\mathbf{x}}$ by

$$S_{\mathbf{x}}(a, b) = \{\mathbf{x}u_{-}(r)g(s), |r| \leq a, |s| \leq b\}.$$

We will refer to $S_{\mathbf{x}}(a, b)$ as a “rectangle”; it isn’t really: it is a quadrilateral bounded by two arcs of geodesic of length $2b$, and by two arcs of negative horocycle, of length respectively ae^b and ae^{-b} (see figure 6).

Further we define the “box” $W_{\mathbf{x}}(a, b, c) \subset \mathbf{X}$ as the region obtained by flowing along positive horocycles from $S_{\mathbf{x}}(a, b)$ until you hit $S_{\mathbf{x}u_{+}(c)}$. Because $S_{\mathbf{x}u_{+}(c)}$ is actually dense in \mathbf{X} , you have to understand the flow as taking place in the universal covering space \mathbf{H} , and then projecting the “box” to \mathbf{X} (see figure 6 again).

Each surface $S_{\mathbf{x}}$ is invariant under geodesic flow, but the “rectangles” $S_{\mathbf{x}}(a, b)$ are not; instead we have

$$S_{\mathbf{x}}(a, b)g(s) = S_{\mathbf{x}g(s)}(e^s a, b).$$

Moreover, the surfaces $S_{\mathbf{x}u_{+}(s)}$ foliate a neighborhood of the positive horocycle $\mathbf{x} * U_{+}$ through \mathbf{x} , and thus there is a function $\alpha_{\mathbf{x}}(\mathbf{y}, t) : \mathbb{R} \rightarrow \mathbb{R}$ for $\mathbf{y} \in S_{\mathbf{x}}$ (for its precise domain, see below) such that

$$\mathbf{y}u_{+}(\alpha_{\mathbf{x}}(\mathbf{y}, t)) \in S_{\mathbf{x}u_{+}(t)}.$$

as sketched in Figure 6.

The function $t \mapsto \alpha_{\mathbf{x}}(\mathbf{y}, t)$ is defined in $[0, T(\mathbf{y}))$ for some $T(\mathbf{y})$ that tends to ∞ as $\mathbf{y} \rightarrow \mathbf{x}$. Moreover, the function is C^{∞} (actually real-analytic) by the implicit function theorem, and $\frac{d}{dt}\alpha_{\mathbf{x}}(\mathbf{y}, t)$ tends to 1 as $\mathbf{y} \rightarrow \mathbf{x}$, so that

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} \frac{\alpha_{\mathbf{x}}(\mathbf{y}, t)}{t} \rightarrow 1.$$

It isn’t often that you can replace the implicit function theorem by an explicit formula, but this does occur here.

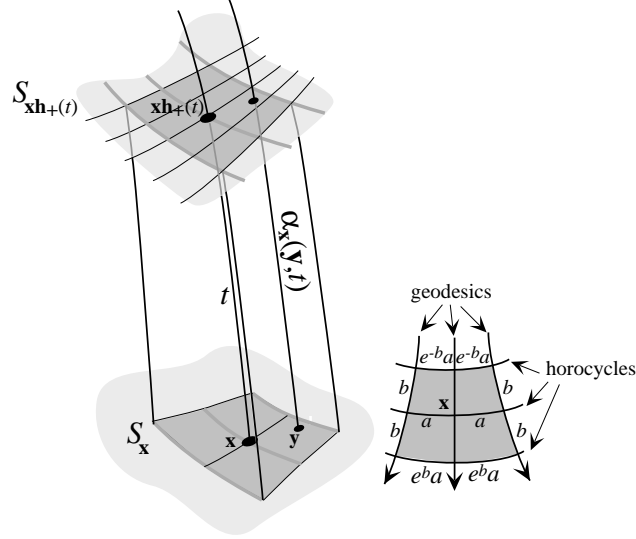


Figure 6: The surfaces $S_{x u_+(s)}$ foliate a neighborhood of the positive horocyclic orbit of \mathbf{x} . Thus, for every t and every $\mathbf{y} \in S_x$ sufficiently close to \mathbf{x} , there is a time $\alpha_x(\mathbf{y}, t)$ such that the positive horocycle $\mathbf{y} u_+(\mathbb{R})$ intersects $S_{x u_+(t)}$.

Lemma 4 *If $\mathbf{y} = \mathbf{x} u_-(r) g(t)$, then*

$$\alpha_x(\mathbf{y}, s) = \frac{s}{e^t(1 - rs)}. \quad (2)$$

In particular, $\alpha_x(\mathbf{y}, s)$ is defined in $\mathbf{y} \in S_x(\frac{\delta}{t}, \delta)$ for all $s < 1/a$ and all b , and

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} \frac{d}{ds} \alpha_x(\mathbf{y}, s) = 1, \quad \lim_{\mathbf{y} \rightarrow \mathbf{x}} \frac{\alpha_x(\mathbf{y}, s)}{s} = 1.$$

Proof. This is a matter of solving the equation

$$\mathbf{x} u_-(r) g(t) u_+(\alpha_x(\mathbf{y}, s)) = \mathbf{x} u_+(s) u_-(\rho) g(\tau),$$

i.e., the matrix equation

$$\begin{bmatrix} 1 & 0 \\ r & 1 \end{bmatrix} \begin{bmatrix} e^{t/2} & 0 \\ 0 & e^{-t/2} \end{bmatrix} \begin{bmatrix} 1 & \alpha \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & s \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \rho & 1 \end{bmatrix} \begin{bmatrix} e^{\tau/2} & 0 \\ 0 & e^{-\tau/2} \end{bmatrix}.$$

This is a system of 3 equations (because the determinants are all 1) in 3 unknowns ρ, τ and α . Just multiply out and check. ■

The central result here is the following:

Lemma 5 *There exists a constant C such that for all $0 < \delta < 1/2$, all $t > 0$, all $\mathbf{y} \in S_{\mathbf{x}}(\frac{\delta}{t}, \delta)$ and all $0 \leq s \leq t$ we have*

$$d(\mathbf{x}u_+(s), \mathbf{y}u_+(\alpha_{\mathbf{x}}(\mathbf{y}, s))) \leq C\delta.$$

Note that $\alpha_{\mathbf{x}}(\mathbf{y}, s)$ is defined for $s \leq t$ when $\mathbf{y} \in S_{\mathbf{x}}(\frac{\delta}{t}, \delta)$ and $\delta \leq 1/2$, since for the factor $1 - rs$ from the denominator of formula 2, we have $r \leq \delta/t$ and $s \leq t$, so $1 - rs \geq 1 - \delta^2 = 3/4$.

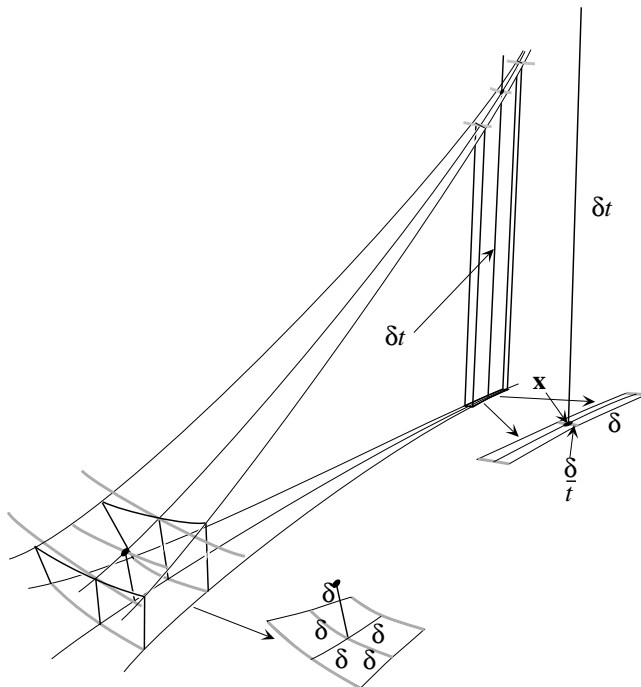


Figure 7: The skinny “rectangle” $S_{\mathbf{x}}(\frac{\delta}{t}, \delta)$ becomes under the geodesic flow for time $\log t$ the “square” $S_{\mathbf{x}g(\log t)}(\delta, \delta)$, and the box $W_{\mathbf{x}}(\frac{\delta}{t}, \delta, \delta t)$ becomes the box $V_{\mathbf{x}}(\delta, \delta, \delta)$. The geometry of $W_{\mathbf{x}}(\delta, \delta, \delta)$ is standard: it depends only on δ . In particular, the positive horocyclic flow from the bottom $S_{\mathbf{x}}(\frac{\delta}{t}, \delta)$ of the box is defined since $\delta \leq 1/2$, hence C^∞ , hence Lipschitz with a universal constant C .

Proof. The proof essentially consists of gazing at Figure 7. Almost everything in that figure comes from the fact that the geodesic flow takes horocycles to horocycles; moreover, geodesic flow for time t maps a segment of

positive horocycle of length l to one of length $e^{-t}l$, and a segment of negative horocycle of length l to one of length $e^t l$.

Two points \mathbf{x}, \mathbf{y} with $\mathbf{y} \in S_{\mathbf{x}}(\delta/t, \delta)$ flow under the geodesic flow for time $\log t$ to two points $\mathbf{x}' = \mathbf{x}g(\log t)$ and $\mathbf{y}' = \mathbf{y}g(\log t)$; note that $\mathbf{y}' \in S_{\mathbf{x}'}(\delta, \delta)$, so certainly $d(\mathbf{x}', \mathbf{y}') \leq 2\delta$. Then under positive horocycle flow (for different times) these points flow to points

$$\mathbf{x}'' = \mathbf{x}'u_+(s\delta) \quad \text{and} \quad \mathbf{y}'' \in S_{\mathbf{x}''}.$$

By the argument in the caption of figure 7, there exists a universal constant C such that $d(\mathbf{x}'', \mathbf{y}'') \leq 2Cd(\mathbf{x}', \mathbf{y}')$. Finally, use the geodesic flow back, i.e., for time $-\log t$, to find points

$$\mathbf{x}''' = \mathbf{x}u_+(s), \mathbf{y}''' = \mathbf{y}u_+(\alpha_{\mathbf{x}}(\mathbf{y}, s)).$$

Since backwards geodesic flow in a single unstable manifold is contracting, we find $d(\mathbf{x}''', \mathbf{y}''') \leq 2C\delta$. ■

6 Geometry of hyperbolic surfaces and cusps

Let X be a complete hyperbolic surface. If such a surface is not compact, it has finitely many *cusps*. Every cusp c is surrounded by closed horocycles, and the open region bounded by the horocycle of length 2 is a neighborhood N_c isometric to the region $\{y \geq 1\}/2\mathbb{Z}$ that is embedded in X , moreover, if c, c' are distinct cusps, then $N_c \cap N_{c'} = \emptyset$.

If X has finite area, then the complement of these neighborhoods is compact:

$$X_c = X - \bigsqcup_{\text{cusps } c \text{ of } X} N_c$$

is a compact set. Denote by \mathbf{X}_c the corresponding part of \mathbf{X} . The injectivity radius is bounded below on \mathbf{X}_c , so there is a number $\delta_{\mathbf{X}} > 0$ such that for every $\mathbf{x} \in \mathbf{X}_c$ the box $W_{\mathbf{x}}(\delta_{\mathbf{X}}, \delta_{\mathbf{X}}, \delta_{\mathbf{X}})$ is embedded.

Now, suppose that \mathbf{X} has finite measure. Then if $(x, \xi) \in \mathbf{X}$ is a point where the positive horocycle is not periodic, the geodesic through (x, ξ) does not go forward to a cusp and hence must enter \mathbf{X}_c infinitely many times.

Let $\mathbf{P}_{\mathbf{X}} \subset \mathbf{X}$ be the set of points defining periodic positive horocycles. Equivalently, $\mathbf{P}_{\mathbf{X}}$ is the union of the stable manifolds of the cusps (for the

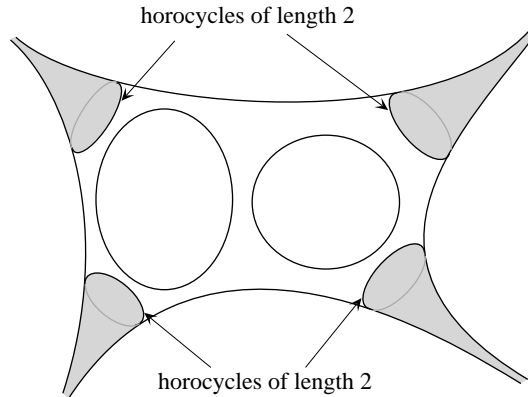


Figure 8: A non-compact complete hyperbolic surface is always non-compact in the same way: it has cusps c with disjoint standard neighborhoods N_c isometric to $\{y \geq 1\}/2\mathbb{Z}$, hence bounded by horocycles of length 2. Note that the only way a geodesic $\gamma(t)$ can stay in such a neighborhood for all $t \geq t_0$ is to head straight to the cusp. Each cusp has a stable manifold in X , and the geodesics that do not return infinitely many times to $X_c = X - \cup_c N_c$ are those that belong to one of these stable manifolds.

geodesic flow). Indeed, if a positive horocycle is closed, it surrounds a cusp, and the geodesic flow from any point of the horocycle goes to the cusp.

Lemma 6 *The set $\mathbf{P}_{\mathbf{X}}$ has measure zero in \mathbf{X} for the measure $\omega_{\mathbf{X}}$.*

Proof. There are finitely many cusps, and each has a stable manifold which is a smooth immersed surface, certainly of 3-dimensional measure 0. ■

7 A sequence of good times

In this section we prove a result, still a bit weaker than theorem 3, though it does prove theorem 3 when X is compact.

Theorem 7 *Let X be a complete hyperbolic surface of finite area, \mathbf{X} be its unit tangent bundle. For all $\mathbf{x} \notin \mathbf{P}_{\mathbf{X}}$, there then exists a sequence $T_n \rightarrow \infty$ such that for any function $f \in C_c(\mathbf{X})$ we have*

$$\lim_{n \rightarrow \infty} \frac{1}{T_n} \int_0^{T_n} f(\mathbf{x}u_+(t)) dt = \int_{\mathbf{X}} f d\omega_{\mathbf{X}}.$$

Choose $\epsilon > 0$, and $f \in C_c(\mathbf{X})$; without loss of generality we may assume $\sup |f| = 1$ and that $\epsilon < 1$.

We have already defined $\delta_{\mathbf{X}}$. We need two more δ 's.

Lemma 8 *There exists $\delta_f > 0$ such that for all $t > 0$, if $\mathbf{z} \in S_{\mathbf{x}}(\delta/t, \delta)$ and $0 \leq s \leq t$, then*

$$|f(\mathbf{x}u_+(s)) - f(\mathbf{z}u_+(\alpha_{\mathbf{x}}(\mathbf{z}, s)))| < \epsilon.$$

Proof. This follows immediately from proposition 5 and the uniform continuity of f . ■

Lemma 9 *There exists δ_{α} such that for all $t > 0$, if $\mathbf{z} \in S_{\mathbf{x}}(\delta/t, \delta)$ and $0 \leq s \leq \delta t$, then*

$$|\alpha'_{\mathbf{x}}(\mathbf{z}, s) - 1| < \epsilon.$$

Proof. One could derive this from the implicit function theorem, but we might as well use our explicit formula (2) for α . For $\mathbf{z} = \mathbf{x}u_-(r)g(u) \in S_{\mathbf{x}}(\delta/t, \delta)$ we have

$$\alpha'_{\mathbf{x}}(\mathbf{y}, s) = \frac{1}{e^u(1-rs)^2},$$

and since $|r| \leq \delta/t$, $|u| \leq \delta$ and $s \leq \delta t$,

$$\frac{1}{e^{\delta}(1+\delta^2)} \leq \alpha'_{\mathbf{x}}(\mathbf{y}, s) \leq \frac{e^{\delta}}{1-\delta^2}.$$

Clearly we can choose δ_{α} so that

$$\left| \frac{1}{e^{\delta_{\alpha}}(1+\delta_{\alpha}^2)} - 1 \right| < \epsilon, \quad \left| \frac{e^{\delta_{\alpha}}}{1-\delta_{\alpha}^2} - 1 \right| < \epsilon.$$

■

Set $\delta = \inf(\delta_{\mathbf{X}}, \delta_f, \delta_{\alpha})$, and $\eta = \omega_{\mathbf{X}}(W_{\mathbf{x}}(\delta, \delta, \epsilon\delta))$.

Proposition 10 *There exists a T_0 and a set $\mathbf{Y} \subset \mathbf{X}$ with $\omega_{\mathbf{X}}(\mathbf{Y}) > 1 - \eta$ such that for all $T > T_0$ and all $\mathbf{y} \in \mathbf{Y}$ we have*

$$\left| \frac{1}{T} \int_0^T f(\mathbf{y}u_+(t))dt - \int_{\mathbf{X}} f d\omega_{\mathbf{X}} \right| < \epsilon.$$

Proof. This follows from the ergodic theorem. This is where we use Theorem 1. ■

We can now choose

1. a sequence $T_n \geq T_0$ tending to infinity such that $\mathbf{x}u_+(T_n) \in \mathbf{X}_c$. This is because the geodesic through \mathbf{x} does not go to a cusp, so it must visit \mathbf{X}_c infinitely many times.
2. a sequence $\mathbf{y}_n \in \mathbf{Y} \cap W_{\mathbf{x}}(\delta/T_n, \delta, \epsilon\delta T_n)$. Indeed, we have

$$\omega_{\mathbf{X}}(W_{\mathbf{x}}(\delta/T_n, \delta, \epsilon\delta T_n)) = \eta,$$

since it is the inverse image of $W_{\mathbf{x}u_+(T_n)}(\delta, \delta, \epsilon\delta)$ by the geodesic flow at time T_n . We have

$$\omega_{\mathbf{X}}(W_{\mathbf{x}u_+(T_n)}(\delta, \delta, \epsilon\delta)) = \eta$$

since $\mathbf{x}u_+(T_n) \in \mathbf{X}_c$ and $\delta \leq \delta_{\mathbf{X}}$. Geodesic flow for a fixed time is a measure-preserving diffeomorphism, so $W_{\mathbf{x}}(\delta/T_n, \delta, \epsilon\delta T_n)$ must intersect \mathbf{Y} which has volume $> 1 - \eta$.

3. sequences $\mathbf{z}_n \in S_{\mathbf{x}}(\delta/T_n, \delta)$ and $\epsilon'_n \leq \epsilon$ such that $\mathbf{z}_n u_+(\epsilon'_n \delta T_n) = \mathbf{y}_n$. This is just what it means to say $\mathbf{y}_n \in W_{\mathbf{x}}(\delta/T_n, \delta, \epsilon\delta T_n)$.

The organizing principle is now to write

$$\begin{aligned} & \left| \frac{1}{T_n} \int_0^{T_n} f(\mathbf{x}u_+(t)) dt - \int_{\mathbf{X}} f(\mathbf{w}) \omega_{\mathbf{X}}(d\mathbf{w}) \right| \leq \\ & \left| \frac{1}{T_n} \int_0^{T_n} f(\mathbf{x}u_+(t)) dt - \frac{1}{T_n} \int_0^{T_n} f(\mathbf{z}_n u_+(\alpha_{\mathbf{x}}(\mathbf{z}_n, t))) dt \right| + \\ & \left| \frac{1}{T_n} \int_0^{T_n} f(\mathbf{z}_n u_+(\alpha_{\mathbf{x}}(\mathbf{z}_n, t))) dt - \frac{1}{T_n} \int_0^{T_n} f(\mathbf{z}_n u_+(\alpha_{\mathbf{x}}(\mathbf{z}_n, t))) \alpha'_{\mathbf{x}}(\mathbf{z}_n, t) dt \right| + \\ & \left| \frac{1}{T_n} \int_0^{\alpha_{\mathbf{x}}(\mathbf{z}_n, T_n)} f(\mathbf{z}_n u_+(s)) ds - \frac{1}{T_n} \int_0^{T_n} f(\mathbf{z}_n u_+(s)) ds \right| + \\ & \left| \frac{1}{T_n} \int_0^{T_n} f(\mathbf{z}_n u_+(s)) ds - \frac{1}{T_n} \int_0^{T_n} f(\mathbf{y}_n u_+(s)) ds \right| + \\ & \left| \frac{1}{T_n} \int_0^{T_n} f(\mathbf{y}_n u_+(s)) ds - \int_{\mathbf{X}} f(\mathbf{w}) \omega_{\mathbf{X}}(d\mathbf{w}) \right|. \end{aligned}$$

To get from the second summand on the right to the third, we use that

$$\int_0^{T_n} f(\mathbf{z}_n u_+(\alpha_{\mathbf{x}}(\mathbf{z}_n, t))) dt = \int_0^{\alpha_{\mathbf{x}}(\mathbf{z}_n, T_n)} f(\mathbf{z}_n u_+(s)) ds$$

by the change of variables formula, setting $s = \alpha_{\mathbf{x}}(\mathbf{z}_n, t)$.

Each of the five terms above need to be bounded in terms of ϵ .

1. Since $\delta < \delta_f$, we have $|f(\mathbf{x}u_+(t)) - f(\mathbf{z}_n u_+(\alpha_{\mathbf{x}}(\mathbf{z}_n, t)))| < \epsilon$, so

$$\left| \frac{1}{T_n} \int_0^{T_n} f(\mathbf{x}u_+(s)) ds - \frac{1}{T_n} \int_0^{T_n} f(\mathbf{z}_n u_+(\alpha_{\mathbf{x}}(\mathbf{z}_n, t))) dt \right| < \epsilon.$$

2. Since $\delta < \delta_\alpha$, we have

$$\begin{aligned} & \left| \frac{1}{T_n} \int_0^{T_n} f(\mathbf{z}_n u_+(\alpha_{\mathbf{x}}(\mathbf{z}_n, t))) dt - \frac{1}{T_n} \int_0^{T_n} f(\mathbf{z}_n u_+(\alpha_{\mathbf{x}}(\mathbf{z}_n, t))) \alpha'_{\mathbf{x}}(\mathbf{z}_n, t) dt \right| \\ & \leq \frac{1}{T_n} \int_0^{T_n} \sup |f| |1 - \alpha'_{\mathbf{x}}(\mathbf{z}_n, t)| dt < \epsilon. \end{aligned}$$

3. From $\delta < \delta_\alpha$, so $|\alpha' - 1| < \epsilon$, we get that $(1 - \epsilon)T_n < \alpha_{\mathbf{x}}(\mathbf{z}_n, T_n) < (1 + \epsilon)T_n$ and hence

$$\left| \frac{1}{T_n} \int_0^{\alpha_{\mathbf{x}}(\mathbf{z}_n, T_n)} f(\mathbf{z}_n u_+(s)) ds - \frac{1}{T_n} \int_0^{T_n} f(\mathbf{z}_n u_+(s)) ds \right| < \epsilon.$$

4. The points \mathbf{z}_n and \mathbf{y}_n are on the same positive horocycle, a distance $\epsilon_n T_n$ apart for some $\epsilon_n \leq \epsilon$. This leads to

$$\begin{aligned} & \left| \frac{1}{T_n} \int_0^{T_n} f(\mathbf{z}_n u_+(s)) ds - \frac{1}{T_n} \int_0^{T_n} f(\mathbf{y}_n u_+(s)) ds \right| \\ & = \left| \frac{1}{T_n} \int_0^{T_n} f(\mathbf{z}_n u_+(s)) ds - \frac{1}{T_n} \int_{\epsilon_n T_n}^{(1+\epsilon_n)T_n} f(\mathbf{z}_n u_+(s)) ds \right| \leq \frac{2\epsilon T_n}{T_n} = 2\epsilon. \end{aligned}$$

5. Since $\mathbf{y}_n \in \mathbf{Y}$ and $T_n > T_0$, we have

$$\left| \frac{1}{T_n} \int_0^{T_n} f(\mathbf{y}_n u_+(s)) ds - \int_{\mathbf{X}} f(\mathbf{w}) \omega_{\mathbf{X}}(d\mathbf{w}) \right| < \epsilon.$$

This ends the proof of theorem 7. \square

8 Proving equidistribution

The sequence T_n was chosen to be times tending to ∞ such that $\mathbf{x}g(T_n) \in \mathbf{X}_c$. Thus if \mathbf{X} is compact, the sequence T_n is an arbitrary sequence tending to infinity, and so equidistribution is proved in that case. Moreover, clearly theorem 7 shows that all non-periodic horocycles are dense in \mathbf{X} . But it doesn't quite prove that they are equidistributed when \mathbf{X} is not compact; perhaps a horocycle could spend an undue amount of time near some cusp, and we could choose a different sequence of times T'_n also tending to infinity which would emphasize the values of f near that cusp. In fact, we will see in Section 9 that something like that does happen for random walks on horocycles.

We will now show that this does not happen for the horocycle flow itself.

Proposition 11 *Let ν be a probability measure on \mathbf{X} invariant under the positive horocycle flow, ergodic for the positive horocycle flow, and such that $\nu(\mathbf{P}_\mathbf{X}) = 0$. Then $\nu = \omega_\mathbf{X}$.*

Proof. Without loss of generality, we can assume that ν is ergodic for the positive horocycle flow since any nonergodic invariant probability measure of this type will be the direct integral of a collection ν_i with $\nu_i(\mathbf{P}_\mathbf{X}) = 0 \forall i$ of ergodic probability measures, so uniqueness for ergodic such measures implies uniqueness for invariant measures of this type. Choose $f \in C_c(\mathbf{X})$, and let $\mathbf{x} \in \mathbf{X}$ be a typical point for ν , i.e., a point of $\mathbf{X} - \mathbf{P}_\mathbf{X}$ such that

$$\int_{\mathbf{X}} f d\nu = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T f(xu_+(s)) ds.$$

By the ergodic theorem, this is true of ν -almost every point, so such points \mathbf{x} certainly exist. Such a point is one for which the horocycle flow is not periodic, so theorem 7 asserts that there exists a sequence $T_n \rightarrow \infty$ such that

$$\int_{\mathbf{X}} f d\nu = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T f(\mathbf{x}u_+(s)) ds = \lim_{n \rightarrow \infty} \frac{1}{T_n} \int_0^{T_n} f(\mathbf{x}u_+(s)) ds = \int_{\mathbf{X}} f d\omega_\mathbf{X}.$$

Since this equality is true for every $f \in C_c(\mathbf{X})$, we have $\nu = \omega_\mathbf{X}$. ■

Now suppose that for some $\mathbf{x} \in \mathbf{X} - \mathbf{P}_\mathbf{X}$ and some $f \in C_c(\mathbf{X})$, we do not have

$$\int_{\mathbf{X}} f d\omega_\mathbf{X} = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T f(\mathbf{x}u_+(s)) ds.$$

We can consider the set of probability measures ν_t defined by

$$\int_{\mathbf{X}} f d\nu_t = \frac{1}{T} \int_0^T f(xu_+(s)) ds.$$

On a non-compact space, the Riesz representation theorem says that set of Borel measures is the dual of the Banach space of the space $C_0(X)$, the space of continuous functions vanishing at ∞ , with the sup norm. The collection of probability measures ν_t is a subset of the unit ball, which is compact for the weak topology. So if $\lim_{t \rightarrow \infty} \nu_t \neq \mu$, there exists a measure $\nu \neq \mu$ and a sequence $t_i \rightarrow \infty$ such that

$$\lim_{i \rightarrow \infty} \nu_{t_i} = \nu$$

in the weak topology.

Clearly ν is invariant under the horocycle flow and ergodic. So it might seem that $\mu \neq \nu$ contradicts proposition 11. There is a difficulty with this argument when \mathbf{X} is not compact. In that case the probability measures do not form a closed subset of the unit ball of $C_0(\mathbf{X})^*$; consider for instance the measures $\delta(x - n)$ on \mathbb{R} ; as $n \rightarrow \infty$ they tend to 0 in the weak topology. Technically, the problem is that we can't evaluate measures on the continuous function 1, since this function doesn't vanish at infinity.

We need to show that ν is a probability measure. This follows from proposition 12 below. For $\rho \leq 2$, let $\mathbf{X}^\rho \subset \mathbf{X}$ be the compact part of \mathbf{X} in which all periodic horocycles have length $\geq \rho$.

Proposition 12 *For any $\epsilon > 0$, there exists $\rho > 0$ such that for all $\mathbf{x} \in \mathbf{X} - \mathbf{P}_{\mathbf{X}}$ we have*

$$\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T \mathbb{1}_{\mathbf{X}^\rho}(xu_+(s)) ds > 1 - \epsilon.$$

Proof. If c is a cusp of X , let N_c^ρ be the neighborhood of c bounded by the horocycle of length ρ . Recall from our discussion of the geometry of hyperbolic surfaces, that N_c^2 is isometric to a standard object: the part of $(2\mathbb{Z}) \backslash \mathbf{H}$ where $y > 1$. Set γ to be the Moebius transformation $\gamma(z) = z/(z+1)$; the standard neighborhood is then isometric to the part of $\langle \gamma \rangle \backslash \mathbf{H}$ where $x^2 + (y-1)^2 \leq 1$. Moreover, any horocycle that doesn't tend to the cusp is equivalent by a change of variable commuting with γ to a horizontal line.

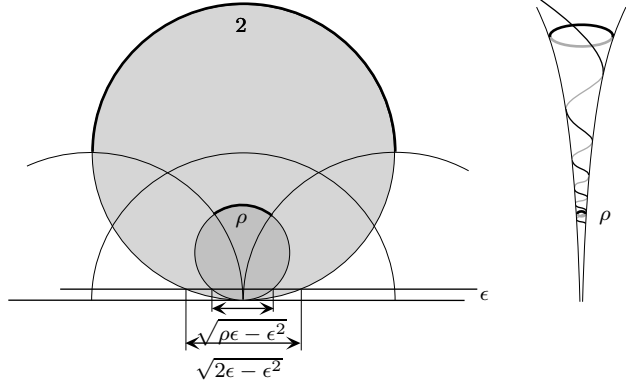


Figure 9: In H a neighborhood of a cusp bounded by a horocycle corresponds to a disc tangent to the x -axis. In the figure on the left, we have represented the cusp by $\langle \gamma \rangle \setminus H$, where without loss of generality we may set $\gamma(z) = z/(1+z)$. Then the disc of radius 1 centered at i corresponds to the neighborhood of the cusp bounded by the horocycle of length 2, and the disc of radius $\rho/2$ centered at $i\rho/2$ corresponds to the neighborhood bounded by a horocycle of length ρ . A horocycle that enters this neighborhood but does not go to the cusp can be, without loss of generality, represented by a line of equation $y = \epsilon$; it goes deeper and deeper into the cusp as $\epsilon \rightarrow 0$. The ratio of times spent in N^ρ to the time spent in $N^2 - N^\rho$ does not become large as the horocycle goes deeper in the cusp, but tends to a ratio depending only on ρ , which tends to 0 as ρ tends to 0. As horocycles go deeper and deeper in the cusp, they spiral more and more tightly in $N^2 - N^\rho$ and still spend approximately the same fraction of time in $N^2 - N^\rho$ as in N^ρ .

Of course lengths on such a horizontal line $y = \epsilon$ depend on ϵ , but ratios of lengths are the same as ratios of euclidean lengths.

A careful look at figure 9 shows that if a horocycle starts in \mathbf{X}_c , goes deep in the cusp, and comes out again, then the ratio of time spent in N^ρ to time spent in $N^2 - N^\rho$ is

$$\frac{\sqrt{\rho\epsilon - \epsilon^2}}{\sqrt{2\epsilon - \epsilon^2} - \sqrt{\rho\epsilon - \epsilon^2}} = \frac{\sqrt{\rho}}{\sqrt{2} - \sqrt{\rho}} + O(\epsilon). \quad (3)$$

Any non-periodic horocycle will eventually enter \mathbf{X}_c ; by taking ρ sufficiently small, we can assure that afterwards it will spend a proportion of its time $< \epsilon$ outside of X^ρ . Proposition 12 follows. ■

Consider the measures

$$\nu_{\mathbf{x},T} = (f \mapsto \frac{1}{T} \int_0^T f(\mathbf{x}u_+(t))dt).$$

Proposition 13 *The accumulation set of $\{\nu_{\mathbf{x},T}, T > 0\}$ consists entirely of probability measures.*

Proof. Every accumulation point μ of the $\nu_{\mathbf{x},T}$ in $C_0(\mathbf{X})^*$ is a measure, and the only thing to show is that $\mu(\mathbf{X}) = 1$. Clearly $\mu(\mathbf{X}) \leq 1$, since for any $f \in C_0(X)$ and any \mathbf{x}, T we have

$$\frac{1}{T} \int_0^T f(\mathbf{x}u_+(t))dt \leq \|f\|_\infty.$$

To see that $\mu(\mathbf{X}) \geq 1$, take $\epsilon > 0$ and ρ as in proposition 12. We can then find a function $f \in C_0(\mathbf{X})$ which coincides with $\mathbb{1}_{\mathbf{X}^\rho}$ on \mathbf{X}^ρ and satisfies $0 \leq f \leq 1$ everywhere. Then

$$\mu(\mathbf{X}) = \sup_{g \in C_0(\mathbf{X})} \frac{\int_{\mathbf{X}} |gd\mu|}{\|g\|_\infty} \tag{4}$$

$$\geq \int_{\mathbf{X}} f d\mu \geq \liminf_{T \rightarrow \infty} \frac{1}{T} \int_0^T f(\mathbf{x}u_+(t))dt \tag{5}$$

$$\geq \liminf_{T \rightarrow \infty} \frac{1}{T} \int_0^T \mathbb{1}_{\mathbf{X}^\rho}(\mathbf{x}u_+(t))dt > 1 - \epsilon. \tag{6}$$

■

There is one last thing to check.

Proposition 14 *A measure μ in the limit set of $\{\nu_{\mathbf{x},T}, T > 0\}$ with $\mathbf{x} \notin \mathbf{P}_{\mathbf{X}}$ satisfies $\mu(\mathbf{P}_{\mathbf{X}}) = 0$.*

Proof. Suppose $\mu(\mathbf{P}_{\mathbf{X}}) > 0$, set $\epsilon = \mu(\mathbf{P}_{\mathbf{X}})/3$ and use proposition 12 to find a corresponding ρ . Find compact subset $\mathbf{Q} \subset \mathbf{P}_{\mathbf{X}}$ with $\mu(\mathbf{Q}) > \frac{2}{3}\mu(\mathbf{P}_{\mathbf{X}})$, and find a time T such that

$$\mathbf{Q}g(T) \cap \mathbf{X}^\rho = \emptyset.$$

This is possible because $\mathbf{P}_{\mathbf{X}}$ consists of points in the stable manifolds of the cusps, so each point can be moved off \mathbf{X}^ρ , and since \mathbf{Q} is compact it will leave \mathbf{X}^ρ under the geodesic flow at some time T .

Let \mathbf{U} be a neighborhood of \mathbf{Q} such that $\mathbf{U}g(T) \cap \mathbf{X}^\rho = \emptyset$. For this neighborhood \mathbf{U} of \mathbf{Q} , as for any neighborhood, there exists a sequence of times $T_n \rightarrow \infty$ such that

$$\frac{\lambda\{t \in [0, T_n] \mid \mathbf{x}u_+(t) \in \mathbf{U}\}}{T_n} > \frac{1}{2}\mu(\mathbf{Q}),$$

where λ is linear measure. Then the horocycle $t \mapsto \mathbf{x}g(T)u_+(t)$ must spend the same proportion of its time in $\mathbf{U}g(T)$, hence outside \mathbf{X}^ρ . But every non-periodic horocycle spends at least a proportion $1 - \mu(\mathbf{P}_X)/3$ in $\mathbf{X} - \mathbf{X}^\rho$, and this is a contradiction. ■

9 Horocycle flow on the modular surface

Let Γ be the 2-congruence subgroup of $\mathrm{SL}_2 \mathbb{R}$, so that $\mathbf{X} = \Gamma \backslash \mathrm{SL}_2 \mathbb{R}$ is the unit tangent bundle over $X = \Gamma \backslash H$, which is the 3-times punctured sphere.

Lemma 15 *The hyperbolic surface X has area 2π , and the subset $X - X^\rho$ has area 3ρ .*

It follows from lemma 15 that for every $\mathbf{x}_0 \notin \mathbf{P}$ there exists for every sequence $\rho_n \rightarrow 0$, for every $\epsilon > 0$ and for every n sufficiently large, a time

$$T_n < (1 + \epsilon) \left(\frac{2\pi}{3\rho_n} \right)$$

such that $\mathbf{x}_0 u_+(T_n) \in X - X^{\rho_n}$.

To use this result, we need to understand the region in H corresponding to X^ρ .

Lemma 16 *The inverse image in H of $X - X^\rho$ is the union of the horodisc $\mathrm{Im} z > 2/\rho$, and the union, for all rational numbers p/q of the discs of radius $\rho/(4q^2)$ tangent to the real axis at p/q .*

So now let us apply the analysis to the horocycle represented by the circle of radius 1 tangent to the real axis at the irrational number α , with $\mathbf{x}_0 = \alpha + 2i$, and set $\rho_n = 1/n$. Further, let us write

$$\mathbf{x}_0 u_+(T) = \left(\alpha - \frac{2T}{T^2 + 1} \right) + i \frac{2}{T^2 + 1}.$$

There thus exists an infinite sequence of rational numbers p_n/q_n and times $T_n < (1 + \epsilon)\frac{2\pi n}{3}$ such that

$$\left| \alpha - \frac{p_n}{q_n} \right| = \frac{2T_n}{T_n^2 + 1} \leq \frac{T_n}{2nq_n^2} < (1 + \epsilon)\frac{\pi n}{3nq_n^2} = (1 + \epsilon)\frac{\pi}{3q_n^2}.$$

This is of course nothing to boast about. It has been known for over 100 years that for every irrational number α , there exist infinitely many coprime numbers p_n, q_n such that

$$\left| \alpha - \frac{p_n}{q_n} \right| < \frac{1}{\sqrt{5}q_n^2},$$

and that $1/\sqrt{5}$ is the smallest number for which this is true [Kin64]. Our analysis only gives the constant $\pi/3$, too large by a factor of more than 2.

One reason to take an interest in this result despite its weakness is that Ratner's theorem has many generalizations to situations where methods leading to the sharp results about diophantine approximations of irrational numbers are not available. In all settings, Ratner's theorem has "diophantine" consequences.

Clearly we cannot do better than improve the constant for all horocycles. But we can use the theory of diophantine approximations to improve the results above for almost every horocycle. In particular we can apply the following theorem.

Theorem 17 [Kin64] *If $g(x) : \mathbb{R}_+^* \rightarrow \mathbb{R}$ is a function such that $g(x)/x$ is increasing, then*

- *there exists infinitely many coprime integers p, q such that*

$$\left| \alpha - \frac{p}{q} \right| < \frac{1}{qg(q)}$$

for almost every α

- *if and only if*

$$\sum_{n=1}^{\infty} \frac{1}{g(n)}$$

is divergent.

Let us see what this says about horocycles. For almost every \mathbf{x}_0 , the horocycle \mathbf{x}_0U_+ lifts to a circle in \mathbf{H} tangent to the real axis at a number α belonging to the set of full measure of theorem 17. Without loss of generality we may assume that $\mathbf{x}_0 = \alpha + 2iR$ for an appropriate R ; it only changes time by an additive constant (and is in any case the worst point to choose on the horocycle).

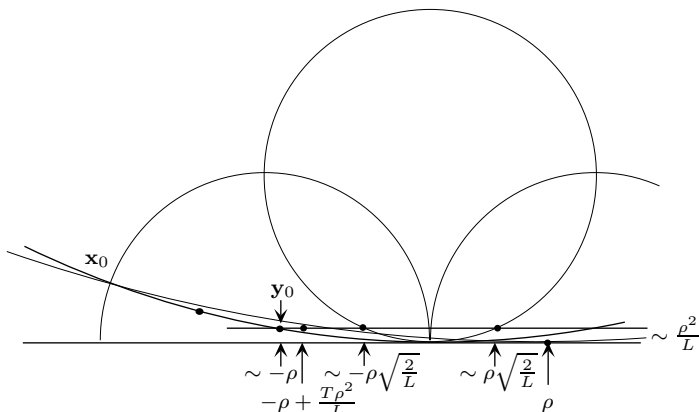


Figure 10: We can lift the horocycle \mathbf{x}_0U_+ to \mathbf{H} ; without loss of generality we may assume that the cusp c is at 0, and that the stabilizer of the cusp is generated by $z \mapsto z/(z+1)$. In that case, the horocycle of length 2 lifts to the circle of radius 1 centered at i , and the horocycle of length L lifts to the circle of radius $L/2$ centered at $Li/2$. We may take \mathbf{x}_0 to be anywhere on this horocycle; it will be convenient to place it at $-2L^2/(L^2+4) + 4Li/(L^2+4)$. In that case one fundamental domain on the horocycle goes from $-2L^2/(L^2+4) + 4Li/(L^2+4)$ to $2L^2/(L^2+4) + 4Li/(L^2+4)$. Our modified horocycle will join x_0 to the point $\rho > 0$ on the real axis. It is much easier to estimate lengths on this horocycle if we send ρ to infinity by a parabolic transformation that fixes 0, and hence all the horocycles tangent to the real axis at 0. If we perform this parabolic transformation, the point x_0 moves to a point \mathbf{y}_0 on its horocycle which is approximately $-\rho + i\rho^2/L$, and the horocycle is a horizontal line, approximately the line $y = \rho^2/L$.

Let p_n/q_n be one of the good approximations to α guaranteed by theorem almostalldiophand let ρ_n be the radius of the negative horocycle surrounding the cusp corresponding to p_n/q_n when the horocycle $\mathbf{x}_0u_+(T)$ is on the vertical line $x = p_n/q_n$, suppose that this occurs for $T = T_n$. Further let us

write

$$\mathbf{x}_0 u_+(T_n) = \xi_n + i\eta_n = 2R \left(\frac{T_n}{T_n^2 + 1} + \frac{i}{T_n^2 + 1} \right).$$

Then we have $T_n = \frac{\xi_n}{\eta_n}$ and $\eta_n = \frac{\rho_n}{2q_n^2}$. It follows that

$$T_n = \frac{\xi_n}{\eta_n} \leq \frac{1/(q_n^2 \log q_n)}{\rho_n/(2q_n^2)} = \frac{2}{\rho_n \log q_n}.$$

Thus on the modular surface, almost every horocycle visits $X - X^{\rho_n}$ at times $T_n \leq \frac{2}{\rho_n \log q_n}$: much earlier than is guaranteed by equidistribution.

This leads to a surprising result due to Breuillard [Bre05]: although non-periodic horocycles are equidistributed, *any* uncentered random walk on the set of non-periodic horocycles almost surely is not.

Theorem 18 *Let μ be a measure on \mathbb{R} with finite expectation and variance:*

$$0 \neq a = \int_{-\infty}^{\infty} t\mu(dt) < \infty \quad \text{and} \quad b^2 = \int_{-\infty}^{\infty} (t - a)^2\mu(dt) < \infty.$$

If $b > 0$, there exists a function $f \in C_c(\mathbf{X})$ with $\int_{\mathbf{X}} f(\mathbf{x})\omega_X(d\mathbf{x}) = 1$ such that for almost every $\mathbf{x}_0 \in \mathbf{X}$ we have

$$\liminf_{m \rightarrow \infty} \int_{-\infty}^{\infty} f(\mathbf{x}_0 u_+(t))\mu^{*m}(dt) = 0.$$

Proof. The measure μ^{*m} is approximately the Gaussian of mean ma and standard deviation \sqrt{mb} . Let us choose m such that $ma = T_n$ for one of the T_n given above, but such that the standard deviation of $\sigma(\mu^{*m}) \sim b\sqrt{m}$ is much smaller than $1/\sqrt{\rho_n}$. This is straightforward, since

$$b\sqrt{m} \sim b\sqrt{\frac{T_n}{a}} \leq \frac{b}{\sqrt{a}}\sqrt{\frac{2}{\rho_n \log q_n}}$$

and $\sigma(\mu^{*m})$ will be much smaller than $1/\sqrt{\rho_n}$ as soon as q_n is large enough.

Recall that it takes time of the order $1/\sqrt{\rho}$ for a horocycle to get from $X - X^\rho$ to X^2 . Thus for the m found above, there are many standard deviations of μ^{*m} around the mean am contained in $X - X^2$. It follows that if $f \in C_c(X)$ satisfies $\int_{\mathbf{X}} f(\mathbf{x})\omega_X(d\mathbf{x}) = 1$ but f has its support in $X - X^2$, we have

$$\liminf_{m \rightarrow \infty} \int_{-\infty}^{\infty} f(\mathbf{x}_0 u_+(t))\mu^{*m}(dt) = 0.$$

This proves that the random walk is not equidistributed. ■

References

- [Bek00] B. Bekka. *Ergodic Theory and Topological Dynamics for Group Actions on Homogeneous Spaces*. Cambridge University Press, 2000.
- [Bre05] E. Breuillard. Local limit theorems and equidistribution of random walks on the heisenberg group. *Geom. Funct. Anal.*, 15(1):35–82, 2005.
- [DS84] S. G. Dani and J. Smillie. Uniform distribution of horocycle orbits for fuchsian groups. *Duke Math. J.*, 51(1):185–194, 1984.
- [Hed36] G. A. Hedlund. Two-dimensional manifolds and transitivity. *Ann. of Math. (2)*, 37(3):534–542, 1936.
- [Hub06] J.H. Hubbard. *Teichmuller Theory and Applications to Geometry, Topology and Dynamics (Vol. 1: Teichmuller Theory)*. Matrix Editions, 2006.
- [Kin64] A. Kinchin. *Continued Fractions*. University of Chicago Press, 1964.
- [Rat92] M. Ratner. Raghunathan’s conjectures for $SL(2, \mathbf{R})$. *Israel J. Math.*, 80(1-2):1–31, 1992.