

CLASSICAL GEOMETRY — LECTURE NOTES

DANNY CALEGARI

1. A CRASH COURSE IN GROUP THEORY

A *group* is an algebraic object which formalizes the mathematical notion which expresses the intuitive idea of *symmetry*. We start with an abstract definition.

Definition 1.1. A *group* is a set G and an operation $m : G \times G \rightarrow G$ called *multiplication* with the following properties:

- (1) m is *associative*. That is, for any $a, b, c \in G$,

$$m(a, m(b, c)) = m(m(a, b), c)$$

and the product can be written unambiguously as abc .

- (2) There is a unique element $e \in G$ called *the identity* with the properties that, for any $a \in G$,

$$ae = ea = a$$

- (3) For any $a \in G$ there is a unique element in G denoted a^{-1} called *the inverse* of a such that

$$aa^{-1} = a^{-1}a = e$$

Given an object with some structural qualities, we can study the symmetries of that object; namely, the set of transformations of the object to itself which preserve the structure in question. Obviously, symmetries can be composed associatively, since the effect of a symmetry on the object doesn't depend on what sequence of symmetries we applied to the object in the past. Moreover, the transformation which does nothing preserves the structure of the object. Finally, symmetries are reversible — performing the opposite of a symmetry is itself a symmetry. Thus, the symmetries of an object (also called the *automorphisms* of an object) are an example of a group.

The power of the abstract idea of a group is that the symmetries can be studied by themselves, without requiring them to be tied to the object they are transforming. So for instance, the same group can act by symmetries of many different objects, or on the same object in many different ways.

Example 1.2. The group with only one element e and multiplication $e \times e = e$ is called the *trivial group*.

Example 1.3. The integers \mathbb{Z} with $m(a, b) = a + b$ is a group, with identity 0.

Example 1.4. The positive real numbers \mathbb{R}^+ with $m(a, b) = ab$ is a group, with identity 1.

Example 1.5. The group with two elements *even* and *odd* and “multiplication” given by the usual rules of addition of even and odd numbers; here *even* is the identity element. This group is denoted $\mathbb{Z}/2\mathbb{Z}$.

Example 1.6. The group of integers mod n is a group with $m(a, b) = a + b \pmod n$ and identity 0. This group is denoted $\mathbb{Z}/n\mathbb{Z}$ and also by C_n , the *cyclic group of length n* .

Definition 1.7. If G and H are groups, one can form the *Cartesian product*, denoted $G \oplus H$. This is a group whose elements are the elements of $G \times H$ where $m : (G \times H) \times (G \times H) \rightarrow G \times H$ is defined by

$$m((g_1, h_1), (g_2, h_2)) = (m_G(g_1, g_2), m_H(h_1, h_2))$$

The identity element is (e_G, e_H) .

Example 1.8. Let S be a regular tetrahedron; label opposite pairs of edges by A, B, C . Then the group of symmetries which preserves the labels is $\mathbb{Z}/2\mathbb{Z} \oplus \mathbb{Z}/2\mathbb{Z}$. It is also known as the *Klein 4-group*.

In all of the examples above, $m(a, b) = m(b, a)$. A group with this property is called *commutative* or *Abelian*. Not all groups are Abelian!

Example 1.9. Let T be an equilateral triangle with sides A, B, C opposite vertices a, b, c in anticlockwise order. The symmetries of T are the reflections in the lines running from the corners to the midpoints of opposite sides, and the rotations. There are three possible rotations, through anticlockwise angles $0, 2\pi/3, 4\pi/3$ which can be thought of as e, ω, ω^2 . Observe that $\omega^{-1} = \omega^2$. Let r_a be a reflection through the line from the vertex a to the midpoint of A . Then $r_a = r_a^{-1}$ and similarly for r_b, r_c . Then $\omega^{-1}r_a\omega = r_c$ but $r_a\omega^{-1}\omega = r_a$ so this group is *not commutative*. It is called the *dihedral group* D_3 and has 6 elements.

Example 1.10. If P is an equilateral n -gon, the symmetries are reflections as above and rotations. This is called the *dihedral group* D_n and has $2n$ elements. The elements are $e, \omega, \omega^2, \dots, \omega^{n-1} = \omega^{-1}$ and r_1, r_2, \dots, r_n where $r_i^2 = e$ for all i , $r_i r_j = \omega^{2(i-j)}$ and $\omega^{-1} r_i \omega = r_{i-1}$.

Example 1.11. The symmetries of an “equilateral ∞ -gon” (i.e. the unique infinite 2-valent tree) defines a group D_∞ , the *infinite dihedral group*.

Example 1.12. The set of 2×2 matrices whose entries are real numbers and whose determinants do not vanish is a group, where multiplication is the usual multiplication of matrices. The set of *all* 2×2 matrices is *not* naturally a group, since some matrices are not invertible.

Example 1.13. The group of permutations of the set $\{1 \dots n\}$ is called the *symmetric group* S_n . A permutation breaks the set up into subsets on which it acts by cycling the members. For example, $(3, 2, 4)(5, 1)$ denotes the element of S_5 which takes $1 \rightarrow 5, 2 \rightarrow 4, 3 \rightarrow 2, 4 \rightarrow 3, 5 \rightarrow 1$. The group S_n has $n!$ elements. A *transposition* is a permutation which interchanges exactly two elements. A permutation is *even* if it can be written as a product of an even number of transpositions, and *odd* otherwise.

Exercise 1.14. Show that the symmetric group is not commutative for $n > 2$. Identify S_3 and S_4 as groups of rigid motions of familiar objects. Show that an even permutation is not an odd permutation, and vice versa.

Definition 1.15. A *subgroup* H of G is a subset such that if $h \in H$ then $h^{-1} \in H$, and if $h_1, h_2 \in H$ then $h_1 h_2 \in H$. With its inherited multiplication operation from G , H is a group. The *right cosets* of H in G are the equivalence classes $[g]$ of elements $g \in G$ where the equivalence relation is given by $g_1 \sim g_2$ if and only if there is an $h \in H$ with $g_1 = g_2 h$.

Exercise 1.16. If H is finite, the number of elements of G in each equivalence class are equal to $|H|$, the number of elements in H . Consequently, if $|G|$ is finite, $|H|$ divides $|G|$.

Exercise 1.17. Show that the subset of even permutations is a subgroup of the symmetric group, known as the alternating group and denoted A_n . Identify A_5 as a group of rigid motions of a familiar object.

Example 1.18. Given a collection of elements $\{g_i\} \subset G$ (not necessarily finite or even countable), the subgroup generated by the g_i is the subgroup whose elements are obtained by multiplying together finitely many of the g_i and their inverses in some order.

Exercise 1.19. Why are only finite multiplications allowed in defining subgroups? Show that a group in which infinite multiplication makes sense is a trivial group. This fact is not as useless as it might seem . . .

Definition 1.20. A group is cyclic if it is generated by a single element. This justifies the notation C_n for $\mathbb{Z}/n\mathbb{Z}$ used before.

Definition 1.21. A homomorphism between groups is a map $f : G_1 \rightarrow G_2$ such that $f(g_1)f(g_2) = f(g_1g_2)$ for any g_1, g_2 in G_1 . The kernel of a homomorphism is the subgroup $K \subset G_1$ defined by $K = f^{-1}(e)$. If $K = e$ then we say f is injective. If every element of G_2 is in the image of f , we say it is surjective. A homomorphism which is injective and surjective is called an isomorphism.

Example 1.22. Every finite group G is isomorphic to a subgroup of S_n where n is the number of elements in G . For, let $b : G \rightarrow \{1, \dots, n\}$ be a bijection, and identify an element g with the permutation which takes $b(h) \rightarrow b(gh)$ for all h .

Definition 1.23. An exact sequence of groups is a (possibly terminating in either direction) sequence

$$\dots \rightarrow G_i \rightarrow G_{i+1} \rightarrow G_{i+2} \rightarrow \dots$$

joined by a sequence of homomorphisms $h_i : G_i \rightarrow G_{i+1}$ such that the image of h_i is equal to the kernel of h_{i+1} for each i .

Definition 1.24. If $a, b \in G$, then bab^{-1} is called the conjugate of a by b , and $aba^{-1}b^{-1}$ is called the commutator of a and b . Abelian groups are characterized by the property that a conjugate of a is equal to a and every commutator is trivial.

Definition 1.25. A subgroup $N \subset G$ is normal, denoted $N \triangleleft G$ if for any $n \in N$ and $g \in G$ we have $gng^{-1} \in N$. A kernel of a homomorphism is normal. Conversely, if N is normal, we can define the quotient group G/N whose elements are equivalence classes $[g]$ of elements in G , and two elements g, h are equivalent iff $g = hn$ for some $n \in N$. The multiplication is given by $m([g], [h]) = [gh]$ and the fact that N is normal says this is well-defined. Thus normal subgroups are exactly kernels of homomorphisms.

Example 1.26. Any subgroup of an abelian group is normal.

Example 1.27. \mathbb{Z} is a normal subgroup of \mathbb{R} . The quotient group \mathbb{R}/\mathbb{Z} is also called the circle group S^1 . Can you see why?

Example 1.28. Let D_n be the dihedral group, and let C_n be the subgroup generated by ω . Then C_n is normal, and $D_n/C_n \cong \mathbb{Z}/2\mathbb{Z}$.

Definition 1.29. If G is a group, the subgroup G_1 generated by the commutators in G is called the commutator subgroup of G . Let G_2 be the subgroup generated by commutators of elements of G with elements of G_1 . We denote $G_1 = [G, G]$ and $G_2 = [G, G_1]$. Define G_i inductively by $G_i = [G, G_{i-1}]$. The elements of G_i are the elements which can be written as products of iterated commutators of length i . If G_i is trivial for some i — that

is, there is some i such that every commutator of length i in G is trivial — we say G is *nilpotent*.

Observe that every G_i is normal, and every quotient G/G_i is nilpotent.

Definition 1.30. If G is a group, let $G^0 = G$ and define $G^i = [G^{i-1}, G^{i-1}]$. If G^i is trivial for some i , we say that G is *solvable*. Again, every G^i is normal and every G/G^i is solvable. Obviously a nilpotent group is solvable.

Definition 1.31. An isomorphism of a group G to itself is called an *automorphism*. The set of automorphisms of G is naturally a group, denoted $\text{Aut}(G)$. There is a homomorphism from $\rho : G \rightarrow \text{Aut}(G)$ where g goes to the automorphism consisting of *conjugation* by g . That is, $\rho(g)(h) = ghg^{-1}$ for any $h \in G$. The automorphisms in the image of ρ are called *inner automorphisms*, and are denoted by $\text{Inn}(G)$. They form a normal subgroup of $\text{Aut}(G)$. The quotient group is called the group of *outer automorphisms* and is denoted by $\text{Out}(G) = \text{Aut}(G)/\text{Inn}(G)$.

Definition 1.32. Suppose we have two groups G, H and a homomorphism $\rho : G \rightarrow \text{Aut}(H)$. Then we can form a new group called *the semi-direct product of G and H* denoted $G \ltimes H$ whose elements are the elements of $G \times H$ and multiplication is given by

$$m((g_1, h_1), (g_2, h_2)) = (g_1g_2, h_1\rho(g_1)(h_2))$$

Observe that H is a normal subgroup of $G \ltimes H$, and there is an exact sequence

$$1 \rightarrow H \rightarrow G \ltimes H \rightarrow G \rightarrow 1$$

Example 1.33. The dihedral group D_n is equal to $\mathbb{Z}/2\mathbb{Z} \ltimes C_n$ where the homomorphism $\rho : \mathbb{Z}/2\mathbb{Z} \rightarrow \text{Aut}(C_n)$ takes the generator of $\mathbb{Z}/2\mathbb{Z}$ to the automorphism $\omega \mapsto \omega^{-1}$, where ω denotes the generator of C_n .

Example 1.34. The group $\mathbb{Z}/2\mathbb{Z} \ltimes \mathbb{R}$ where the nontrivial element of $\mathbb{Z}/2\mathbb{Z}$ acts on \mathbb{R} by $x \mapsto -x$ is isomorphic to the group of *isometries* (i.e. 1–1 and distance preserving transformations) of the real line. It contains D_∞ as a subgroup.

Exercise 1.35. Find an action of $\mathbb{Z}/2\mathbb{Z}$ on the group S^1 so that D_n is a subgroup of $\mathbb{Z}/2\mathbb{Z} \ltimes S^1$ for every n .

Example 1.36. The group whose elements consist of words in the alphabet a, b, A, B subject to the equivalence relation that when one of aA, Aa, bB, Bb appear in a word, they may be removed, so for example

$$aBaAbb \sim aBbb \sim ab$$

A word in which none of these special subwords appears is called *reduced*; it is clear that the equivalence classes are in 1–1 correspondence with reduced words. Multiplication is given by concatenation of words. The identity is the empty word, $A = a^{-1}, B = b^{-1}$. In general, the inverse of a word is obtained by reversing the order of the letters and changing the case. This is called the *free group F_2 on two generators*, in this case the letters a, b . It is easy to generalize to the *free group F_n on n generators*, given by words in letters a_1, \dots, a_n and their “inverse letters” A_1, \dots, A_n . One can also denote the letters A_i by the “letters” a_i^{-1} .

Exercise 1.37. Let G be an arbitrary group and $g_1, g_2 \dots g_n$ a finite subset of G . Show that there is a unique homomorphism from $F_n \rightarrow G$ sending $a_i \rightarrow g_i$.

Example 1.38. If we have an alphabet consisting of letters a_1, \dots, a_n and their inverses, we can consider a collection of words in these letters r_1, \dots, r_m . If R denotes the subgroup of F_n generated by the r_i and all their conjugates, then R is a normal subgroup of F_n and we can form the quotient F_n/R . This is denoted by

$$\langle a_1, \dots, a_n \mid r_1, \dots, r_m \rangle$$

and an equivalent description is that it is the group whose elements are words in the a_i and their inverses modulo the equivalence relation that two words are equivalent if they are equivalent in the free group, or if one can be obtained from the other by inserting or deleting some r_i or its inverse as a subword somewhere. The a_i are the *generators* and the r_i the *relations*. Groups defined this way are very important in topology. Notice that a *presentation* of a group in terms of generators and relations is far from unique.

Definition 1.39. A group G is *finitely generated* if there is a finite subset of G which generates G . This is equivalent to the property that there is a surjective homomorphism from some F_n to G . A group G is *finitely presented* if it can be expressed as $\langle A \mid R \rangle$ for some finite set of generators A and relations R .

Exercise 1.40. Let G be any finite group. Show that G is finitely presented.

Exercise 1.41. Let F_2 be the free group on generators x, y . Let $i : F_2 \rightarrow \mathbb{Z}$ be the homomorphism which takes $x \rightarrow 1$ and $y \rightarrow 1$. Show that the kernel of i is not finitely generated.

Exercise 1.42. (Harder). Let $i : F_2 \oplus F_2 \rightarrow \mathbb{Z}$ be the homomorphism which restricts on either factor to i in the previous exercise. Show that the kernel of i is finitely generated but not finitely presented.

Definition 1.43. Given groups G, H the *free product of G and H* , denoted $G * H$, is the group of words whose letters alternate between elements of G and H , with concatenation as multiplication, and the obvious proviso that the identity is in either G or H . It is the unique group with the *universal property* that there are injective homomorphisms $i_G : G \rightarrow G * H$ and $i_H : H \rightarrow G * H$, and given any other group I and homomorphisms $j_G : G \rightarrow I$ and $j_H : H \rightarrow I$ there is a *unique* homomorphism c from $G * H$ to I satisfying $c \circ i_G = j_G$ and $c \circ i_H = j_H$.

Exercise 1.44. Show that $*$ defines an associative and commutative product on groups up to isomorphism, and

$$F_n = \mathbb{Z} * \mathbb{Z} * \dots * \mathbb{Z}$$

where we take n copies of \mathbb{Z} in the product above.

Exercise 1.45. Show that $\mathbb{Z}/2\mathbb{Z} * \mathbb{Z}/2\mathbb{Z} \cong D_\infty$.

Remark 1.46. Actually, one can extend $*$ to *infinite* (even uncountable) products of groups by the universal property. If one has an arbitrary set S the *free group generated by S* is the free product of a collection of copies of \mathbb{Z} , one for each element of S .

Exercise 1.47. (Hard). Every subgroup of a free group is free.

Definition 1.48. A *topological group* is a group which is also a space (i.e. we understand what continuous maps of the space are) such that $m : G \times G \rightarrow G$ and $i : G \rightarrow G$, the multiplication and inverse maps respectively, are *continuous*. If G is a *smooth manifold* (see appendix for definition) and the maps m and i are smooth maps, then G is called a *Lie group*.

Remark 1.49. Actually, the usual definition of Lie group requires that G be a *real analytic manifold* and that the maps m and i be real analytic. A real analytic manifold is like a smooth manifold, except that the co-ordinate transformations between charts are required to be real analytic, rather than merely smooth. It turns out that any *connected, locally connected, locally compact* (see appendix for definition) topological group is actually a Lie group.

2. MODEL GEOMETRIES IN DIMENSION TWO

2.1. The Euclidean plane.

2.1.1. Euclid's axioms.

Notation 2.1. The Euclidean plane will be denoted by \mathbb{E}^2 .

Euclid, who taught at Alexandria in Egypt and lived from about 325 BC to 265 BC, is thought to have written 13 famous mathematical books called the *Elements*. In these are found the earliest (?) historical example of the *axiomatic method*. Euclid proposed 5 postulates or axioms of geometry, from which all true statements about the Euclidean plane were supposed to inevitably follow. These axioms were as follows:

- (1) A straight line segment can be drawn joining any two points.
- (2) Any straight line segment is contained in a unique straight line.
- (3) Given any straight line segment, a circle can be drawn having the segment as radius and one endpoint as center.
- (4) All right angles are congruent.
- (5) One and only one line can be drawn through a point parallel to a given line.

The terms *point*, *line*, *plane* are supposed to be primitive concepts, in the sense that they can't be described in terms of simpler concepts. Since they are not defined, one is not supposed to use one's personal notions or intuitions about these objects to prove theorems about them; one strategy to achieve this end is to replace the terms by other terms (Hilbert's suggestion is *glass*, *beer mat*, *table*; Queneau's is *word*, *sentence*, *paragraph*) or even nonsense terms. The point is not that intuition is worthless (it is *not*), but that by proving theorems about objects by only using the properties expressed in a list of axioms, the proof immediately applies to any other objects which satisfy the same list of axioms, including collections of objects that one might not have originally had in mind. In this way, our ordinary geometric intuitions of space and movement can be used to reason about objects far from our immediate experience. One important remark to make is that, by modern standards, Euclid's foundations are far from rigorous. For instance, it is implicit in the statement of the axioms that *angles can be added*, but nowhere is it said what properties this addition satisfies; angles are *not* numbers, neither are lengths, but they have properties in common with them.

2.1.2. *A closer look at the fourth postulate.* Notice that Euclid does not define "congruence". A working definition is that two figures X and Y in a space Z are *congruent* if there is a transformation of Z which takes X to Y . But which transformations are allowed? By including certain kinds of transformations and excluding others, we can drastically affect the flavor of the geometry in question. If not enough transformations are allowed, distinct objects are incomparable and one cannot say anything meaningful about them. If too many transformations are allowed, differences collapse and the supply of distinct objects to investigate dries up. One way of reformulating the fourth postulate is to say that space is

homogeneous: that is, the properties of an object do not depend on where it is placed in space. Most of the spaces we will encounter in the sequel will be homogeneous.

2.1.3. *A closer look at the parallel postulate.* The fifth axiom above is also known as the *parallel postulate*. To decode it, one needs a workable definition of parallel. The “usual” definition is that two distinct lines are parallel if and only if they do not intersect. So the postulate says that given a line l and a point p disjoint from l , there is a unique line l_p through p such that l_p and l are disjoint. Historically, this axiom was seen as unsatisfying, and much effort was put into attempts to show that it followed inevitably as a consequence of the other four axioms. Such an attempt was doomed to failure, for the simple reason that there are interpretations of the “undefined concepts” point, line, plane which satisfy the first four axioms but which do *not* satisfy the fifth. If we say that given l and p there is *no* line l_p through p which does not intersect l , we get *elliptic geometry*. If we say that given l and p there are *infinitely many* lines l_p through p which do not intersect l , we get *hyperbolic geometry*. Together with Euclidean geometry, these geometries will be the main focus of this course.

2.1.4. *Symmetries of \mathbb{E}^2 .* What are the “allowable” transformations in Euclidean geometry? That is, what are the transformations of \mathbb{E}^2 which preserve the geometrical properties which characterize it? These special transformations are called the *symmetries* (also called *automorphisms*) of \mathbb{E}^2 ; they form a *group*, which we will denote by $\text{Aut}(\mathbb{E}^2)$. A symmetry of \mathbb{E}^2 takes lines to lines, and preserves angles, but a symmetry of \mathbb{E}^2 does *not* have to preserve lengths. A symmetry can either preserve or reverse orientation. Basic symmetries include *translations*, *rotations*, *reflections*, *dilations*. It turns out that all symmetries of \mathbb{E}^2 can be expressed as simple combinations of these.

Exercise 2.2. *Let $f : \mathbb{E}^2 \rightarrow \mathbb{E}^2$ be orientation-reversing. Show that there is a unique line l such that f can be written as $g \circ r$ where r is a reflection in l and g is an orientation-preserving symmetry which fixes l , in which case g is either a translation parallel to l or a dilation whose center is on l . A reflection in l followed by a translation parallel to l is also called a glide reflection.*

Denote by $\text{Aut}^+(\mathbb{E}^2)$ the orientation-preserving symmetries, and by $\text{Isom}^+(\mathbb{E}^2)$ the orientation-preserving symmetries which are also distance-preserving.

Exercise 2.3. *Suppose $f : \mathbb{E}^2 \rightarrow \mathbb{E}^2$ is in $\text{Aut}^+(\mathbb{E}^2)$ but not in $\text{Isom}^+(\mathbb{E}^2)$. Then there is a unique point p fixed by f , and we can write f as $r \circ d$ where d is a dilation with center p and r is a rotation with center p .*

Exercise 2.4. *Suppose $f \in \text{Isom}^+(\mathbb{E}^2)$. Then either f is a rotation or a translation, and it is a translation exactly when it does not have a fixed point. In either case, f can be written as $r_1 \circ r_2$ where r_i is a reflection in some line l_i . f is a translation exactly when l_1 and l_2 are parallel.*

These exercises show that any distance-preserving symmetry can be written as a product of at most 3 reflections. An interesting feature of these exercises is that they can be established *without using the parallel postulate*. So they describe true facts (where relevant) about elliptic and about hyperbolic geometry. So, for instance, a distance preserving symmetry of the hyperbolic plane can be written as a product $r_1 \circ r_2$ of reflections in lines l_1, l_2 , and this transformation has a fixed point if and only if the lines l_1, l_2 intersect.

Exercise 2.5. *Verify that the group of orientation-preserving similarities of \mathbb{E}^2 which fix the origin is isomorphic to \mathbb{C}^* , the group of non-zero complex numbers with multiplication*

as the group operation. Verify too that the group of translations of \mathbb{E}^2 is isomorphic to \mathbb{C} with addition as the group operation.

Exercise 2.6. Verify that the group $\text{Aut}^+(\mathbb{E}^2)$ of orientation-preserving similarities of \mathbb{E}^2 is isomorphic to $\mathbb{C}^* \times \mathbb{C}$ where \mathbb{C}^* acts on \mathbb{C} by multiplication. In this way identify $\text{Aut}^+(\mathbb{E}^2)$ with the group of 2×2 complex matrices of the form

$$\begin{bmatrix} \alpha & \beta \\ 0 & 1 \end{bmatrix}$$

and $\text{Isom}^+(\mathbb{E}^2)$ with the subgroup where $|\alpha| = 1$.

2.2. The 2–sphere.

2.2.1. Elliptic geometry.

Notation 2.7. The 2–sphere will be denoted by \mathbb{S}^2 .

A very interesting “re–interpretation” of Euclid’s first 4 axioms gives us elliptic geometry. A *point* in elliptic geometry consists of *two antipodal points* in \mathbb{S}^2 . A *line* in elliptic geometry consists of a *great circle* in \mathbb{S}^2 . The *antipodal* map $i : \mathbb{S}^2 \rightarrow \mathbb{S}^2$ is the map which takes any point to its antipodal point. A “line” or “point” with the interpretation above is invariant (as a set) under i , so we may think of the action as all taking place in the “quotient space” \mathbb{S}^2/i . An object in this quotient space is just an object in \mathbb{S}^2 which is invariant as a set by i . Any two great circles intersect in a pair of antipodal points, which is a single “point” in \mathbb{S}^2/i . If we think of \mathbb{S}^2 as a subset of \mathbb{E}^3 , a great circle is the intersection of the sphere with a plane in \mathbb{E}^3 through the origin. A pair of antipodal points is the intersection of the sphere with a line in \mathbb{E}^3 through the origin. Thus, the geometry of \mathbb{S}^2/i is equivalent to the geometry of planes and lines in \mathbb{E}^3 . A plane in \mathbb{E}^3 through the origin is perpendicular to a unique line in \mathbb{E}^3 through the origin, and vice-versa. This defines a “duality” between lines and points in \mathbb{S}^2/i ; so for any theorem one proves about lines and points in elliptic geometry, there is an analogous “dual” theorem with the idea of “line” and “point” interchanged. Let d denote the transformation which takes points to lines and vice versa.

Circles and angles make sense on a sphere, and one sees that the first 4 axioms of Euclid are satisfied in this model.

As distinct from Euclidean geometry where there are symmetries which change lengths, there is a natural length scale on the sphere. We set the diameter equal to 2π .

2.2.2. *Spherical trigonometry.* An example of this duality (and a justification of the choice of length scale) is given by the following

Lemma 2.8 (Spherical law of sines). *If T is a spherical triangle with side-lengths A, B, C and opposite angles α, β, γ , then*

$$\frac{\sin(A)}{\sin(\alpha)} = \frac{\sin(B)}{\sin(\beta)} = \frac{\sin(C)}{\sin(\gamma)}$$

Notice that the triangle $d(T)$ has side lengths $(\pi - \alpha), (\pi - \beta), (\pi - \gamma)$ and angles $(\pi - A), (\pi - B), (\pi - C)$. Notice too that $\sin(t) \approx t$ for small t , so that if T is a very small triangle, this formula approximates the sine rule for Euclidean space. Let \mathbb{S}_t^2 denote the sphere scaled to have diameter $2\pi t$; then the term $\frac{\sin(A)}{\sin(\alpha)}$ in the spherical sine rule should be replaced with $\frac{t \sin(t^{-1}A)}{\sin(\alpha)}$. In this way we may think of \mathbb{E}^2 as the “limit” as $t \rightarrow \infty$ of \mathbb{S}^2 .

Exercise 2.9. *Prove the spherical law of sines. Think of the sides of T as the intersection of \mathbb{S}^2 with planes π_i through the origin in \mathbb{E}^3 , intersecting in lines l_i in \mathbb{E}^3 . Then the lengths A, B, C are the angles between the l_i and the angles α, β, γ are the angles between the planes π_i .*

2.2.3. *The area of a spherical triangle.* If L is a lune of \mathbb{S}^2 between the longitude 0 and the longitude α , then the area of L is 2α .

Now, let T be an arbitrary spherical triangle. If T is bounded by sides l_i which meet at vertices v_i then we can extend the sides l_i to great circles which cut up \mathbb{S}^2 into eight regions. Each pair of lines bound two lunes, and the six lunes so produced fall into two sets of three which intersect exactly along the triangle T and the antipodal triangle $i(T)$. It follows that we can calculate the area of \mathbb{S} as follows

$$4\pi = \text{area}(\mathbb{S}^2) = \sum \text{area}(\text{lunes}) - 4 \text{area}(T) = 4(\alpha + \beta + \gamma) - 4 \text{area}(T)$$

In particular, we have the beautiful formula, which is a special case of the Gauss–Bonnet theorem:

Theorem 2.10. *Let T be a spherical triangle with angles α, β, γ . Then*

$$\text{area}(T) = \alpha + \beta + \gamma - \pi$$

Notice that as T gets very small and the area $\rightarrow 0$, the sum of the angles of T approach π . Thus in the limit, we have Euclidean geometry in which the sum of the angles of a triangle are π . The angle formula for Euclidean triangles is equivalent to the parallel postulate.

Exercise 2.11. *Derive a formula for the area of a spherical polygon with n vertices in terms of the angles.*

Exercise 2.12. *Using the spherical law of sines and the area formula, calculate the area of a regular spherical n -gon with sides of length t .*

2.2.4. *Kissing numbers — the Newton–Gregory problem.* How many balls of radius 1 can be arranged in \mathbb{E}^3 so that they all touch a fixed ball of radius 1? It is understood that the balls are non-overlapping, but they may touch each other at a single point; figuratively, one says that the balls are “kissing” or “osculating” (from the Latin word for kiss), and that one wants to know the *kissing number* in 3-dimensions.

Exercise 2.13. *What is the kissing number in 2-dimensions? That is, how many disks of radius 1 can be arranged in \mathbb{E}^2 so that they all touch a fixed disk of radius 1?*

This question first arose in a conversation between Isaac Newton and David Gregory in 1694. Newton thought 12 balls was the maximum; Gregory thought 13 might be possible. It is quite easy to arrange 12 balls which all touch a fixed ball — arrange the centers at the vertices of a regular icosahedron. If the distance from the center of the icosahedron to the vertices is 2, it turns out the distance between adjacent vertices is ≈ 2.103 , so this configuration can be physically realized (i.e. there is no overlapping). The problem is that there is some slack in this configuration — the balls roll around, and it is unclear whether by packing them more tightly there would be room for another ball.

Suppose we have a configuration of non-overlapping spheres S_i all touching the central sphere S . Let v_i be the points on S where they all touch. The non-overlapping condition is exactly equivalent to the condition that no two of the v_i are a distance of less than $\frac{\pi}{3}$ apart. If some of the S_i are loose, roll them around on the surface until they come into contact with other S_j ; it's clear that we can roll “loose” S_i around until every S_i touches *at least*

two other S_j, S_k . If S_i touches S_{j_1}, \dots, S_{j_n} then join v_i to v_{j_1}, \dots, v_{j_n} by segments of great circles on S . This gives a decomposition of S into spherical polygons, every edge of which has length $\frac{\pi}{3}$. It's easy to see that no polygon has 6 or more sides (why?).

Let f_n be the number of faces with n sides. Then there are $\frac{3}{2}f_3 + 2f_4 + \frac{5}{2}f_5$ edges, since every edge is contained in two faces. Recall Euler's formula for a polygonal decomposition of a sphere

$$\text{faces} - \text{edges} + \text{vertices} = 2$$

so the number of vertices is $2 + \frac{1}{2}f_3 + f_4 + \frac{3}{2}f_5$

Exercise 2.14. Show that the largest spherical quadrilateral or pentagon with side lengths $\frac{\pi}{3}$ is the regular one. Use your formula for the area of such a polygon and the fact above to show that the kissing number is 12 in 3-dimensions. This was first shown in the 19th century.

Exercise 2.15. Show what we have implicitly assumed: namely that a connected nonempty graph in S^2 with embedded edges, and no vertices of valence 1, has polygonal complementary regions.

Remark 2.16. In 1951 Schutte and van der Waerden ([8]) found an arrangement of 13 unit spheres which touches a central sphere of radius $r \approx 1.04556$ where r is a root of the polynomial

$$4096x^{16} - 18432x^{12} + 24576x^{10} - 13952x^8 + 4096x^6 - 608x^4 + 32x^2 + 1$$

This r is thought to be optimal.

2.2.5. *Reflections, rotations, involutions; $SO(3)$.* By thinking of \mathbb{S}^2 as the unit sphere in \mathbb{E}^3 , and by thinking of points and lines in \mathbb{S}^2 as the intersection of the sphere with lines and planes in \mathbb{E}^3 we see that symmetries of \mathbb{S}^2 extend to linear maps of \mathbb{E}^3 to itself which fix the origin. These are expressed as 3×3 matrices. The condition that a matrix M induce a symmetry of \mathbb{S}^2 is exactly that it preserves distances on \mathbb{S}^2 ; equivalently, it preserves the angles between lines through the origin in \mathbb{E}^3 . Consequently, it takes orthonormal frames to orthonormal frames. (A *frame* is another word for a *basis*.)

Any frame can be expressed as a 3×3 matrix F , where the columns give each of the vectors. F is orthonormal if $F^t F = \text{id}$. If M preserves orthonormality, then $F^t M^t M F = \text{id}$ for every orthonormal F ; in particular, $M^t M = F F^t = \text{id}$. Observe that each of these transformations actually induces a symmetry of \mathbb{S}^2 ; in particular, we can identify the set of symmetries of \mathbb{S}^2 with the set of orthonormal frames in \mathbb{E}^3 , which can be identified with the set of 3×3 matrices M satisfying $M^t M = \text{id}$. It is easy to see that such matrices form a *group*, known as the *orthogonal group* and denoted $O(3)$. The *subgroup* of orientation-preserving matrices (those with determinant 1) are denoted $SO(3)$ and called the *special orthogonal group*.

Exercise 2.17. Show that every element of $O(3)$ has an eigenvector with eigenvalue 1 or -1 . Deduce that a symmetry of \mathbb{S}^2 is either a rotation, a reflection, or a product $s \circ r$ where r is reflection in some great circle l and s is a rotation which fixes that circle. (How is this like a "glide reflection"?) In particular, every symmetry of \mathbb{S}^2 is a product of at most three reflections. Compare with the Euclidean case.

2.2.6. *Algebraic groups.* Once we have "algebraized" the geometry of \mathbb{S}^2 by comparing it with the group of matrices $O(3)$ we can generalize in unexpected ways. Let \mathbb{A} denote the *field* of real algebraic numbers. That is, the elements of \mathbb{A} are the real roots of polynomials with rational coefficients. If $a, b \in \mathbb{A}$ and $b \neq 0$ then $a + b, a - b, ab, a/b$ are all in \mathbb{A} (this

is the defining property of a field). There is a natural subgroup of $O(3)$ denoted $O(3, \mathbb{A})$ called the *3-dimensional orthogonal group over \mathbb{A}* which consists of the 3×3 matrices M with entries in \mathbb{A} satisfying $M^t M = \text{id}$. Observe that this, too is a group. If $p = (0, 0, 1)$ we can consider the subset $S^2(\mathbb{A})$, the set of points in \mathbb{S}^2 which are translates of p by elements of $O(3, \mathbb{A})$.

We think of $S^2(\mathbb{A})$ as the points in a funny kind of space. Let

$$S^1(\mathbb{A}) = S^2(\mathbb{A}) \cap \{z = 0\}$$

and define the set of *lines* in $S^2(\mathbb{A})$ as the translates of $S^1(\mathbb{A})$ by elements of $O(3, \mathbb{A})$.

First observe that if q, r are any two points in $S^2(\mathbb{A})$ thought of as vectors in \mathbb{E}^3 then the length of their vector cross-product is in \mathbb{A} . If q, r are two points in $S^2(\mathbb{A})$ then together with 0 they lie on a plane $\pi(q, r)$. Then the triple

$$\left(q, q \times \frac{q \times r}{\|q \times r\|}, \frac{q \times r}{\|q \times r\|} \right)$$

is an orthonormal frame with co-ordinates in \mathbb{A} . If we think of this triple as an element of $O(3, \mathbb{A})$ then the image of $S^1(\mathbb{A})$ contains q and r . Thus there is a “line” in $S^2(\mathbb{A})$ through q and r . (Here \times denotes the usual cross product of vectors.)

Exercise 2.18. Show that the set $S^2(\mathbb{A})$ is exactly the set of points in $\mathbb{S}^2 \subset \mathbb{E}^3$ with co-ordinates in \mathbb{A} .

Exercise 2.19. Explore the extent to which Euclid’s axioms hold or fail to hold for $S^2(\mathbb{A})$ or $S^2(\mathbb{A})/i$. What if one replaces \mathbb{A} with another field, like \mathbb{Q} ?

Let $O(2, \mathbb{A})$ be the subgroup of $O(3, \mathbb{A})$ which fixes the vector $p = (0, 0, 1)$. Then if $M \in O(2, \mathbb{A})$ and $N \in O(3, \mathbb{A})$, $N(p) = NM(p)$. So we can identify $S^2(\mathbb{A})$ with the quotient space $O(3, \mathbb{A})/O(2, \mathbb{A})$, which is the set of equivalence classes $[N]$ where $N \in O(3, \mathbb{A})$ and $[N] \sim [N']$ if and only if there is a $M \in O(2, \mathbb{A})$ with $N = N'M$. In general, if F denotes an arbitrary field, we can think of the group $O(3, F)$ as the set of 3×3 matrices with entries in F such that $M^t M = \text{id}$. This contains $O(2, F)$ naturally as the subgroup which fixes the vector $(0, 0, 1)$, and we can study the quotient space $O(3, F)/O(2, F)$ as a geometrical space in its own right. Notice that $O(3, F)$ acts by symmetries on this space, by $M \cdot [N] \rightarrow [MN]$. The quotient space is called a *homogeneous space of $O(3, F)$* .

Exercise 2.20. Let F be the field of integers modulo multiples of 2. What is the group $O(3, F)$? How many points are in the space $O(3, F)/O(2, F)$?

In general, if we have a group of matrices defined by some algebraic condition, for instance $\det(M) = 1$ or $M^t M = \text{id}$ or $M^t J M = J$ for $J = \begin{bmatrix} 0 & I \\ -I & 0 \end{bmatrix}$ etc. then we can consider the group of matrices satisfying the condition with coefficients in some field. This is called an *algebraic group*. Many properties of certain algebraic groups are independent of the coefficient field. An algebraic group over a finite field is a finite group; such finite groups are very important, and form the building blocks of “most” of finite group theory.

2.2.7. *Quaternions and the group \mathbb{S}^3* . Recall that the *quaternions* are elements of the 4-dimensional real vector space spanned by $1, i, j, k$ with multiplication which is linear in each factor, and on the basis elements is given by

$$ij = k, jk = i, ki = j, i^2 = j^2 = k^2 = -1$$

This multiplication is associative. The *norm* of a quaternion, denoted by

$$\|a_1 + a_2i + a_3j + a_4k\| = (a_1^2 + a_2^2 + a_3^2 + a_4^2)^{1/2}$$

is equal to the length of the corresponding vector (a_1, a_2, a_3, a_4) in \mathbb{R}^4 . Norms are multiplicative. That is, $\|\alpha\beta\| = \|\alpha\|\|\beta\|$. The non-zero quaternions form a *group* under multiplication; the unit quaternions, which correspond exactly to the unit length vectors in \mathbb{R}^4 , are a subgroup which is denoted \mathbb{S}^3 . Let π be the set of quaternions of the form $1 + ai + bj + ck$. Then π is a copy of \mathbb{R}^3 , and is the *tangent space* to the sphere \mathbb{S}^3 of unit norm quaternions at 1. The group \mathbb{S} acts on π by

$$\alpha \cdot z = \alpha^{-1}z\alpha$$

for $z \in \pi$. Since it preserves lengths, the image is isomorphic to a subgroup of $SO(3, \mathbb{R})$, which can be thought of as the group of orthogonal transformations of π . In fact, this homomorphism is *surjective*. Moreover, the kernel is exactly the center of \mathbb{S}^3 , which is ± 1 . That is, we have the isomorphism $\mathbb{S}^3 / \pm 1 \cong SO(3, \mathbb{R})$.

Exercise 2.21. Write down the formula for an explicit homomorphism, in terms of standard quaternionic co-ordinates for \mathbb{S}^3 and matrix co-ordinates for $SO(3, \mathbb{R})$.

In general, the conjugation action of a Lie group on its tangent space at the identity is called the *adjoint action* of the group. Since the group of linear transformations of this vector space is a matrix group, this gives a homomorphism of the Lie group to a matrix group.

2.3. The hyperbolic plane.

2.3.1. The problem of models.

Notation 2.22. The hyperbolic plane will be denoted by \mathbb{H}^2 .

The sphere is relatively easy to understand and visualize because there is a very nice model of it in Euclidean space: the unit sphere in \mathbb{E}^3 . Symmetries of the sphere extend to symmetries of the ambient space, and distances and angles in the sphere are what one expects from the ambient embedding. No such model exists of the hyperbolic plane. Bits of the hyperbolic plane can be isometrically (i.e. in a distance-preserving way) embedded, but not in such a way that the natural symmetries of the plane can be realized as symmetries of the embedding. However, if we are willing to look at embeddings which distort distances, there are some very nice models of the hyperbolic plane which one can play with and get a good feel for.

2.3.2. *The Poincaré Model:* Suppose we imagine the world as being circumscribed by the unit circle in the plane. In order to prevent people from falling off the edge, we make the edges very cold. As everyone knows, objects shrink when they get cold, so people wandering around on the disk would get smaller and smaller as they approached the edge, so that its apparent distance (to them) would get larger and larger and they could never reach it. Technically, the “length elements” at the point (x, y) are

$$\left(\frac{2dx}{(1-x^2-y^2)}, \frac{2dy}{(1-x^2-y^2)} \right)$$

or in polar co-ordinates, the “length elements” at the point r, θ are

$$\left(\frac{2dr}{(1-r^2)}, \frac{2rd\theta}{(1-r^2)} \right)$$

This is called the *Poincaré metric* on the unit disk, and the disk with this metric is called the *Poincaré model of the hyperbolic plane*. With this choice of metric, the length of a radial line from the origin to the point $(r_1, 0)$ is

$$\int_0^{r_1} \frac{2dr}{(1-r^2)} = \log\left(\frac{1+r}{1-r}\right)\Bigg|_0^{r_1} = \log\left(\frac{1+r_1}{1-r_1}\right)$$

If γ is any other path from the origin to $(r_1, 0)$ whose Euclidean length is l , then its length in the hyperbolic metric is

$$\text{hyperbolic length of } \gamma = \int_0^l \frac{2dt}{(1-r^2(\gamma(t)))}$$

where $r(\gamma(t))$ is the distance from the point $\gamma(t)$ to the origin. Obviously, $l \geq r_1$ and $r(\gamma(t)) \leq t$ with equality if and only if γ is a Euclidean straight line. This implies that the shortest curve in the hyperbolic metric from the origin to a point in the disk is the Euclidean straight line.

Notice that for a point p at Euclidean distance ϵ from the boundary circle, the ratio of the hyperbolic to the Euclidean metric is

$$\frac{2}{1-(1-\epsilon)^2} \approx \frac{1}{\epsilon}$$

for sufficiently small ϵ .

Exercise 2.23. (Hard). Let E be a simply-connected (i.e. without holes) domain in \mathbb{R}^2 bounded by a smooth curve γ . Define a “metric” on E as follows. Let f be a smooth, nowhere zero function on E which is equal to $\frac{1}{\text{dist}(p,\gamma)}$ for all p sufficiently close to γ . Let the length elements on E be given by $(dx f, dy f)$ where (dx, dy) are the usual Euclidean length elements. Show that there is a continuous, 1–1 map $\phi : E \rightarrow D$ which distorts the lengths of curves by a bounded amount. That is, there is a constant $K > 0$ such that for any curve α in E ,

$$\frac{1}{K} \text{length}_D(\phi(\alpha)) \leq \text{length}_E(\alpha) \leq K \text{length}_D(\phi(\alpha))$$

Definition 2.24. The circle ∂D is called the *circle at infinity* of D , and is denoted S_∞^1 . A point in S_∞^1 is called an *ideal point*.

Now think of the unit disk D as the set of complex numbers of norm ≤ 1 . Let α, β be two complex numbers with $|\alpha|^2 - |\beta|^2 = 1$. The set of matrices of the form

$$M = \begin{bmatrix} \alpha & \bar{\beta} \\ \beta & \bar{\alpha} \end{bmatrix}$$

form a group, called the *special unitary group* $SU(1, 1)$. This is exactly the group of complex linear transformations of \mathbb{C}^2 which preserves the function $v(z, w) = |z|^2 - |w|^2$ and have determinant 1. These are the matrices M of determinant 1 satisfying $\bar{M}^t J M = J$

where $J = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$.

Exercise 2.25. Define $U(1, 1)$ to be the group of 2×2 complex matrices M with $\bar{M}^t J M = J$ with J as above, and no condition on the determinant. Find the most general form of a matrix in $U(1, 1)$. Show that these are exactly the matrices whose column vectors are an orthonormal basis for the “norm” defined by v .

Now, there is a natural action of $SU(1, 1)$ on D by

$$\begin{bmatrix} \alpha & \bar{\beta} \\ \beta & \bar{\alpha} \end{bmatrix} \cdot z \rightarrow \frac{\alpha z + \bar{\beta}}{\beta z + \bar{\alpha}}$$

Observe that two matrices which differ by ± 1 act on D in the same way. So the action descends to the quotient group $SU(1, 1)/\pm 1$ which is denoted by $PSU(1, 1)$ for the *projective special unitary group*.

Observe that this transformation preserves the boundary circle. Furthermore it takes lines and circles to lines and circles, and preserves angles of intersection between them; in particular, it permutes segments of lines and circles perpendicular to ∂D .

Exercise 2.26. *Show that there is a transformation in $PSU(1, 1)$ taking any point in the interior of D to any other point.*

Exercise 2.27. *Show that the subgroup of $PSU(1, 1)$ fixing any point is isomorphic to a circle. Deduce that we can identify D with the coset space $PSU(1, 1)/S^1$; i.e. D is a homogeneous space for $PSU(1, 1)$.*

Exercise 2.28. *Show that the action of $PSU(1, 1)$ preserves the Poincaré metric in D^2 . Deduce that the shortest hyperbolic path between any two points is through an arc of a circle orthogonal to ∂D or, if the points are on the same diameter, by a segment of this diameter.*

2.3.3. *The upper half-space model.* The *upper half-space*, denoted \mathbb{H} , is the set of points $x, y \in \mathbb{R}^2$ with $y > 0$. Suppose now that the real line is chilled, so that distances in this model are scaled in proportion to the distance to the boundary. That is, in (x, y) coordinates the “length elements” of the metric are

$$\left(\frac{dx}{y}, \frac{dy}{y} \right)$$

Observe that translations parallel to the x -axis preserve the metric, and are therefore isometries. Also, dilations centered at points on the x -axis preserve the metric too. The length of a vertical line segment from (x, y_1) to (x, y_2) is

$$\int_{y_1}^{y_2} \frac{dy}{y} = \log \frac{y_2}{y_1}$$

A similar argument to before shows that this is the shortest path between these two points.

The group of 2×2 matrices with real coefficients and determinant 1 is called the *special linear group*, and denoted $SL(2, \mathbb{R})$ or $SL(2)$ if the coefficients are understood. These matrices act on the upper half-plane, thought of as a domain in \mathbb{C} , by

$$\begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix} \cdot z \rightarrow \frac{\alpha z + \beta}{\gamma z + \delta}$$

Again, two matrices which differ by a constant multiple act on \mathbb{H} in the same way, so the action descends to the quotient group $SL(2, \mathbb{R})/\pm 1$ which is denoted by $PSL(2, \mathbb{R})$ for the *projective special linear group*.

As before, these transformations take lines and circles to lines and circles, and preserve the real line.

Exercise 2.29. *Show that the action of $PSL(2, \mathbb{R})$ preserves the metric on \mathbb{H} . Deduce that the shortest hyperbolic path between any two points in the upper half-space is through an arc of a circle orthogonal to \mathbb{R} or, if the points are on the same vertical line, through a segment of this line.*

Exercise 2.30. Find a transformation from D to \mathbb{H} which takes the Poincaré metric on the disk to the hyperbolic metric on \mathbb{H} . Deduce that these models describe “the same” geometry. Find an explicit isomorphism $PSL(2, \mathbb{R}) \cong PSU(1, 1)$.

2.3.4. *The hyperboloid model.* In \mathbb{R}^3 let H denote the sheet of the hyperboloid $x^2 + y^2 - z^2 = -1$ with z positive. Let $O(2, 1)$ denote the set of 3×3 matrices with real entries which preserve the function $v(x, y, z) = x^2 + y^2 - z^2$, and $SO(2, 1)$ the subgroup with determinant 1. Equivalently, $O(2, 1)$ is the group of real matrices M such that $M^t J M = J$ where

$$J = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -1 \end{bmatrix}$$

Then $SO(2, 1)$ preserves the sheet H .

Definition 2.31. A vector v in \mathbb{R}^3 is *timelike* if $v^t J v < 0$, *spacelike* if $v^t J v > 0$ and *lightlike* if $v^t J v = 0$. The *Lorentz length* of a vector is $(v^t J v)^{1/2}$, denoted $\|v\|$, and can be positive, zero, or imaginary. The *timelike angle* between two timelike vectors v, w is

$$\eta(v, w) = \cosh^{-1} \left(\frac{v^t J w}{\|v\| \|w\|} \right)$$

Compare this with the usual angle between two vectors in \mathbb{R}^3 :

$$\nu(v, w) = \cos^{-1} \left(\frac{v^t w}{\|v\| \|w\|} \right)$$

where in this equation $\|\cdot\|$ denotes the usual length of a vector. Notice that H is exactly the set of vectors of Lorentz length i , just as \mathbb{S}^2 is the set of vectors of usual length 1 (this has led some people to comment that the hyperbolic plane should be thought of as a “sphere of imaginary radius”). Since distances between points in \mathbb{S}^2 are defined as the angle between the vectors, it makes sense to define distances in H as the timelike angle between vectors. For two vectors $v, w \in H$ the formula above simplifies to

$$\eta(v, w) = \cosh^{-1}(-v^t J w)$$

Exercise 2.32. Let K be the group of matrices of the form

$$\begin{bmatrix} \cos(\alpha) & \sin(\alpha) & 0 \\ -\sin(\alpha) & \cos(\alpha) & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

and A the group of matrices of the form

$$\begin{bmatrix} \cosh(\gamma) & 0 & \sinh(\gamma) \\ 0 & 1 & 0 \\ \sinh(\gamma) & 0 & \cosh(\gamma) \end{bmatrix}$$

Show that every element of $SO(2, 1)$ can be expressed as $k_1 a k_2$ for some $k_1, k_2 \in K$ and $a \in A$; that is, we can write $SO(2, 1) = K A K$. How unique is such an expression?

Notice the group K above is precisely the stabilizer of the point $(0, 0, 1) \in H$. Thus we can identify H with the homogeneous space $SO(2, 1)/K$.

Exercise 2.33. Let K' be the subgroup of $PSL(2, \mathbb{R})$ consisting of matrices of the form $\begin{bmatrix} \cos(\alpha) & \sin(\alpha) \\ -\sin(\alpha) & \cos(\alpha) \end{bmatrix}$ and A' the subgroup of matrices of the form $\begin{bmatrix} s & 0 \\ 0 & s^{-1} \end{bmatrix}$. Find an

isomorphism from $SO(2, 1)$ to $PSL(2, \mathbb{R})$ taking K to K' and A to A' . (Careful! The isomorphism $K \rightarrow K'$ might not be the one you first think of. . .)

Remark 2.34. The isomorphisms $PSU(1, 1) \cong PSL(2, \mathbb{R}) \cong SO(2, 1)$ are known as *exceptional isomorphisms*, and one should not assume that the matter is so simple in higher dimensions. Such exceptional isomorphisms are rare and are a very powerful tool, since difficult problems about one of the groups can become simpler when translated into a problem about another of the groups.

After identifying $PSL(2, \mathbb{R})$ with $SO(2, 1)$ we can identify their homogeneous spaces $PSL(2, \mathbb{R})/K'$ and $SO(2, 1)/K$. This identification of H with \mathbb{H} shows that H is an equivalent model of the hyperbolic plane. The straight lines in H are the intersection of planes in \mathbb{R}^3 through the origin with H . In many ways, the hyperboloid model of the hyperbolic plane is the closest to the model of \mathbb{S}^2 as the unit sphere in \mathbb{R}^3 .

Exercise 2.35. Show that the identification of H with \mathbb{H} preserves metrics.

Two vectors $v, w \in \mathbb{R}^3$ are *Lorentz orthogonal* if $v^t J w = 0$. It is easy to see that if v is timelike, any orthogonal vector w is spacelike. If $v \in H$ and w is a tangent vector to H at v , then $v^t J w = 0$, since the derivative $\frac{d}{dt} \|v + tw\|$ should be equal to 0 at $t = 0$ (by the definition of a tangent vector). For two spacelike vectors v, w which span a spacelike vector space, the value of $\frac{v^t J w}{\|v\| \|w\|} < 1$, so $\eta(v, w)$ is an imaginary number. The *spacelike angle* between v and w is defined to be $-i\eta(v, w)$. It is a fact that the hyperbolic angle between two tangent vectors at a point in H is exactly equal to their spacelike angle. The “proof” of this fact is just that the symmetries of the space H preserve this spacelike angle, and the total spacelike angle of a circle is 2π . Since these two properties uniquely characterize hyperbolic angles, the two notions of angle agree.

The fact that hyperbolic lengths and angles can be expressed so easily in terms of trigonometric functions and linear algebra makes the hyperboloid model the model of choice for doing hyperbolic trigonometry.

2.3.5. The Klein (projective) model. Let D_1 be the disk consisting of points in \mathbb{R}^3 with $x^2 + y^2 \leq 1$ and $z = 1$. Then we can project H to D_1 along rays in \mathbb{R}^3 which pass through the origin.

Exercise 2.36. Verify that this stereographic projection is 1–1 and onto.

This projection takes the straight lines in H to (Euclidean) straight lines in D_1 . This gives us a new model of the hyperbolic plane as the unit disk, whose points are usual points, and whose lines are exactly the segments of Euclidean lines which intersect D_1 . One should be wary that Euclidean angles in this model do not accurately depict the true hyperbolic angles. In this model, two lines l_1, l_2 are perpendicular under the following circumstances:

- If l_1 passes through the origin, l_2 is perpendicular to l_1 if and only if it is perpendicular in the Euclidean sense.
- Otherwise, let m_1, m_2 be the two tangent lines to ∂D_1 which pass through the endpoints of l_1 . Then l_1 and l_2 are perpendicular if and only if l_2, m_1, m_2 intersect in a point.

Exercise 2.37. Verify that l_1 is perpendicular to l_2 if and only if l_2 is perpendicular to l_1 .

It is easy to verify in this model that all Euclid’s axioms but the fifth are satisfied.

The relationship between the Klein model and the Poincaré model is as follows: we can map the Poincaré disk to the northern hemisphere of the unit sphere by stereographic projection. This preserves angles and takes lines and circles to lines and circles. In this model, the straight lines are exactly the arcs of circles perpendicular to the equator. Then the Klein model and this (curvy) Poincaré model are related by placing thinking of the Klein disk as the flat Euclidean disk spanning the equator, and mapping D_1 to the upper hemisphere by projecting points along lines parallel to the z -axis. This takes lines in D_1 to the intersection of the upper hemisphere with vertical planes. These intersections are the circular arcs which are perpendicular to the equator, so this map takes straight lines to straight lines as it should.

Exercise 2.38. Write down the metric in D_1 for the Klein model. Using this formula, calculate the hyperbolic distance between the center and a point at radius r_1 in D_1 .

Exercise 2.39. Show that Pappus' theorem is true in the hyperbolic plane; this theorem says that if a_1, a_2, a_3 and b_1, b_2, b_3 are points in two lines l_1 and l_2 , then the six line segments joining the a_i to the b_j for $i \neq j$ intersect in three points which are collinear.

2.3.6. *Hyperbolic trigonometry.* Hyperbolic and spherical geometry are two sides of the same coin. For many theorems in spherical geometry, there is an analogous theorem in hyperbolic geometry. For instance, we have the

Lemma 2.40 (Hyperbolic law of sines). If T is a hyperbolic triangle with sides of length A, B, C opposite angles α, β, γ then

$$\frac{\sinh(A)}{\sin(\alpha)} = \frac{\sinh(B)}{\sin(\beta)} = \frac{\sinh(C)}{\sin(\gamma)}$$

Exercise 2.41. Prove the hyperbolic law of sines by using the hyperboloid model and trying to imitate the vector proof of the spherical law of sines.

2.3.7. *The area of a hyperbolic triangle.* The parallels between spherical and hyperbolic geometry are carried further by the theorem for the area of a hyperbolic triangle. We relax slightly the notion of a triangle: we allow some or all of the vertices of our triangle to be ideal points. If all three vertices are ideal, we say that we have an *ideal triangle*. Notice that since every hyperbolic straight line is perpendicular to S_∞^1 , the angle of a triangle at an ideal point is 0.

Theorem 2.42. Let T be a hyperbolic triangle with angles α, β, γ . Then

$$\text{area}(T) = \pi - \alpha - \beta - \gamma$$

Proof: In the upper half-space model, let T be the triangle with one ideal point at ∞ and two ordinary points at $(\cos(\alpha), \sin(\alpha))$ and $(\cos(\pi - \beta), \sin(\pi - \beta))$ in Euclidean coordinates, where α, β are both $\leq \pi/2$. Such a triangle has angles $0, \alpha, \beta$. The hyperbolic area is

$$\begin{aligned} & \int_{x=\cos(\pi-\beta)}^{\cos(\alpha)} \left(\int_{y=1-x^2}^{\infty} \frac{1}{y^2} dy \right) dx \\ &= \int_{x=\cos(\pi-\beta)}^{\cos(\alpha)} \frac{1}{1-x^2} dx \\ &= -\cos^{-1}(x) \Big|_{\cos(\pi-\beta)}^{\cos(\alpha)} = \pi - \alpha - \beta \end{aligned}$$

In particular, a triangle with one or two ideal points satisfies the formula. Now, for an arbitrary triangle T with angles α, β, γ we can dissect an ideal triangles with all angles 0 into T and three triangles, each of which has two ideal points, and whose third angle is one of $\pi - \alpha, \pi - \beta, \pi - \gamma$. That is,

$$\text{area}(T) = \pi - (\pi - (\pi - \alpha)) - (\pi - (\pi - \beta)) - (\pi - (\pi - \gamma)) = \pi - \alpha - \beta - \gamma$$

□

2.3.8. Projective geometry. The group $PSL(2, \mathbb{R})$ acts in a natural way on another space called the *projective line*, denoted \mathbb{RP}^1 . This is the space whose points are the lines through the origin in \mathbb{R}^2 . Equivalently, this is the quotient of the space $\mathbb{R}^2 - 0$ by the equivalence relation that $(x, y) \sim (\lambda x, \lambda y)$ for any $\lambda \in \mathbb{R}^*$. The unit circle maps 2–1 to \mathbb{RP}^1 , so one sees that \mathbb{RP}^1 is itself a circle. The natural action of $PSL(2, \mathbb{R})$ on \mathbb{RP}^1 is the *projectivization* of the natural action of $SL(2, \mathbb{R})$ on \mathbb{R}^2 . That is,

$$\begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} \alpha x + \beta y \\ \gamma x + \delta y \end{bmatrix}$$

We can write an equivalence class (x, y) unambiguously as x/y , where we write ∞ when $y = 0$; in this way, we can naturally identify \mathbb{RP}^1 with $\mathbb{R} \cup \infty$. One sees that in this formulation, this is exactly the action of $PSL(2, \mathbb{R})$ on the ideal boundary of \mathbb{H}^2 in the upper half–space model. That is, *the geometry of \mathbb{RP}^1 is hyperbolic geometry at infinity*. Observe that for any two triples of points a_1, a_2, a_3 and b_1, b_2, b_3 in \mathbb{RP}^1 which are circularly ordered, there is a *unique* element of $PSL(2, \mathbb{R})$ taking a_i to b_i . For a *four–tuple* of points a_1, a_2, a_3, a_4 let γ be the transformation taking a_1, a_2, a_3 to $0, 1, \infty$. Then $\gamma(a_4)$ is an invariant of the 4–tuple, called the *cross–ratio* of the four points, denoted $[a_1, a_2, a_3, a_4]$. Explicitly,

$$[a_1, a_2, a_3, a_4] = \frac{(a_1 - a_3)(a_2 - a_4)}{(a_1 - a_2)(a_3 - a_4)}$$

The subgroup of $PSL(2, \mathbb{R})$ which fixes a point in $\mathbb{R} \cup \infty$ is isomorphic to the group of *orientation preserving similarities of \mathbb{R}* , which we could denote by $\text{Aut}^+(\mathbb{R})$. This group is isomorphic to $\mathbb{R}^+ \ltimes \mathbb{R}$, where \mathbb{R}^+ acts on \mathbb{R} by multiplication. We can think of \mathbb{RP}^1 as the homogeneous space $PSL(2, \mathbb{R})/\mathbb{R}^+ \ltimes \mathbb{R}$.

Projective geometry is the geometry of *perspective*. Imagine that we have a transparent glass pane, and we are trying to capture a landscape by setting up the pane and painting the scenery on the pane as it appears to us. We could move the pane to the right or left; this would translate the scene left or right respectively. We could move the pane closer or further away; this would shrink or magnify the image. Or we could rotate the pane and ourselves so that the sun doesn't get in our eyes. The horizon in our picture is \mathbb{RP}^1 , and the transformations we can perform on the image is precisely the projective group $PSL(2, \mathbb{R})$.

2.3.9. Elliptic, parabolic, hyperbolic isometries. There are three different kinds of transformations in $PSL(2, \mathbb{R})$ which can be distinguished by their action on S_∞^1 .

Definition 2.43. A non–trivial element $\alpha \in PSL(2, \mathbb{R})$ is *elliptic, parabolic* or *hyperbolic* if it has respectively 0, 1 or 2 fixed points in S_∞^1 . These cases can be distinguished by the property that $|tr(\gamma)|$ is $<, =$ or > 2 , where tr denotes the trace of a matrix representative of γ .

An elliptic transformation has a unique fixed point in \mathbb{H}^2 and acts as a rotation about that point. A hyperbolic transformation fixes the geodesic running between its two ideal

fixed points and acts as a translation along this geodesic. Furthermore, the points in \mathbb{H}^2 moved the shortest distance by the transformation are exactly the points on this geodesic.

A parabolic transformation has no analogue in Euclidean or Spherical geometry. It has no fixed point, but moves points an arbitrarily short amount. In some sense, it is like a “rotation about an ideal point”. Two elliptic elements are conjugate iff they rotate about their respective fixed points by the same amount. This angle of rotation is equal to $\cos^{-1}(|\text{tr}(\gamma)/2|)$. Two hyperbolic elements are conjugate iff they translate along their geodesic by the same amount. This translation length is equal to $\cosh^{-1}(|\text{tr}(\gamma)/2|)$. If a, b are any two parabolic elements then either a and b are conjugate or a^{-1} and b are conjugate.

Exercise 2.44. *Prove the claims made in the previous paragraph.*

Exercise 2.45. *Recall the subgroups K' and A' defined in the prequel. Let N denote the group of matrices of the form $\begin{bmatrix} 1 & t \\ 0 & 1 \end{bmatrix}$. Show that every element of $PSL(2, \mathbb{R})$ can be expressed as kan for some $k \in K'$, $a \in A'$ and $n \in N$. How unique is this expression? This is an example of what is known as the KAN or Iwasawa decomposition.*

Exercise 2.46. *Consider the group $SL(n, \mathbb{R})$ of $n \times n$ matrices with real entries and determinant 1. Let K be the subgroup $SO(n, \mathbb{R})$ of real $n \times n$ matrices M satisfying $M^t M = id$. Let A be the subgroup of diagonal matrices. Let N be the subgroup of matrices with 1's on the diagonal and 0's below the diagonal. Show A is abelian and N is nilpotent. Further, K is compact (see appendix), thought of as a topological subspace of the space \mathbb{R}^{n^2} of $n \times n$ matrices. Show that there is a KAN decomposition for $SL(n, \mathbb{R})$.*

2.3.10. *Horocircular geometry.* In the Poincaré disk model, the Euclidean circles in D which are tangent to S^1_∞ are special; they are called *horocircles* and one can think of them as *circles of infinite radius*. If the point of tangency is taken to be ∞ in the upper half-space model, these circles correspond to the horizontal lines in the upper half-plane. Observe that a parabolic element fixes the family of horocircles tangent to its fixed ideal point, and acts on each of them by translation.

3. TESSELLATIONS

3.1. The topology of surfaces.

3.1.1. *Gluing polygons.* Certain computer games get around the constraint of a finite screen by means of a trick: when a spaceship comes to the left side of the screen, it disappears and “reappears” on the right side of the screen. Likewise, an asteroid which disappears beyond the top of the screen might reappear menacingly from the bottom. The screen can be represented by a square whose sides are *labelled in pairs*: the left and right sides get one label, the top and bottom sides get another label. These labels are instructions for obtaining an idealized topological space from the flat screen: the left and right sides can be glued together to make a cylinder, then the top and bottom sides can be glued together to make a torus (the surface of a donut). Actually, we have to be somewhat careful: there are *two* ways to glue two sides together; an unambiguous instruction must specify the orientation of each edge.

If we have a collection of polygons P_i to be glued together along pairs of edges, we can imagine a graph Γ whose vertices are the polygons, and whose edges are the pairs of edges in the collection. We can glue in any order. If we first glue along the edges corresponding to a *maximal tree* in Γ , the result of this first round of gluing will produce a connected

polygon. Thus, without loss of generality, it suffices to consider gluings of the sides of a single polygon.

Definition 3.1. A *surface* is a 2–dimensional manifold. That is, a Hausdorff topological space with a countable basis, such that every point has a neighborhood homeomorphic to the open unit disk in \mathbb{R}^2 (see appendix for definitions). A *piecewise–linear surface* is a surface obtained from a countable collection of polygons by glueing together the edges in pairs, in such a way that only finitely many edges are incident to any vertex.

Exercise 3.2. *Why is the finiteness condition imposed on vertices?*

The following theorem was proved by T. Rado in 1924 (see [6]):

Theorem 3.3. *Any surface is homeomorphic to a piecewise–linear surface. Any compact surface is homeomorphic to a piecewise–linear surface made from only finitely many polygons.*

Definition 3.4. A surface is *oriented* if there is an unambiguous choice of “top” and “bottom” side of each polygon which is compatible with the glueing. i.e. an orientable surface is “two–sided”.

3.1.2. *The fundamental group.* The definition of the fundamental group of a surface Σ requires a choice of a *basepoint* in Σ . Let $p \in \Sigma$ be such a point.

Definition 3.5. Define $\Omega_1(\Sigma, p)$ to be the space of continuous maps $c : S^1 \rightarrow \Sigma$ sending $0 \in S^1$ to $p \in \Sigma$ (here we think of S^1 as $I/0 \sim 1$). Two such maps c_1, c_2 are called *homotopic* if there is a map $C : S^1 \times I \rightarrow \Sigma$ such that

- (1) $C(\cdot, 0) = c_1(\cdot)$.
- (2) $C(\cdot, 1) = c_2(\cdot)$.
- (3) $C(0, \cdot) = p$.

Exercise 3.6. *Show that the relation of being homotopic is an equivalence relation on $\Omega_1(\Sigma, p)$.*

In fact, $\Omega_1(\Sigma, p)$ has the natural structure of a topological space; with respect to this topological structure, the equivalence classes determined by the homotopy relation are the path–connected components.

The importance of the relation of homotopy equivalence is that the equivalence classes form a *group*:

Definition 3.7. The *fundamental group* of Σ with basepoint p , denoted $\pi_1(\Sigma, p)$, has as elements the equivalence classes

$$\pi_1(\Sigma, p) = \Omega_1(\Sigma, p)/\text{homotopy equivalence}$$

with the group operation defined by

$$[c_1] \cdot [c_2] = [c_1 * c_2]$$

where $c_1 * c_2$ denotes the map $S^1 \rightarrow \Sigma$ defined by

$$c_1 * c_2(t) = \begin{cases} c_1(2t) & \text{for } t \leq 1/2 \\ c_2(2t - 1) & \text{for } t \geq 1/2 \end{cases}$$

where the identity is given by the equivalence class $[e]$ of the constant map $e : S^1 \rightarrow p$, and inverse is defined by $[c]^{-1} = [i(c)]$, where $i(c)$ is the map defined by $i(c)(t) = c(1 - t)$.

Exercise 3.8. Check that $[c][c]^{-1} = [e]$ with the definitions given above, so that $\pi_1(\Sigma, p)$ really is a group.

For Σ a piecewise-linear surface with basepoint v a vertex of Σ , define $O_1(\Sigma)$ to be the space of polygonal loops γ from v to v contained in the edges of Σ . Such a loop consists of a sequence of oriented edges

$$\gamma = e_1, e_2, \dots, e_n$$

where e_1 starts at v and e_n ends there, and e_i ends where e_{i+1} starts. Let $r(e)$ denote the same edge e with the opposite orientation. For a polygonal path γ (not necessarily starting and ending at v) let γ^{-1} denote the path obtained by reversing the order and the orientation of the edges in γ .

One can perform an elementary move on a polygonal loop γ , which is one of the following two operations:

- If there is a vertex w which is the endpoint of some e_i , and α is any polygonal path beginning at w , we can insert or delete $\gamma\alpha^{-1}$ between e_i and e_{i+1} . That is,

$$e_1, e_2, \dots, e_i, \gamma, \alpha^{-1}, e_{i+1}, \dots, e_n \longleftrightarrow e_1, e_2, \dots, e_n$$

- If γ is a loop which is the boundary of a polygonal region, and starts and ends at a vertex w which is the endpoint of some e_i , then we can insert or delete γ between e_i and e_{i+1} . That is,

$$e_1, e_2, \dots, e_i, \gamma, e_{i+1}, \dots, e_n \longleftrightarrow e_1, e_2, \dots, e_n$$

Definition 3.9. Define the combinatorial fundamental group of Σ , denoted $p_1(\Sigma, v)$, to be the group whose elements are the equivalence classes

$$O_1(\Sigma, v)/\text{elementary moves}$$

With the group operation defined by $[\alpha][\beta] = [\alpha\beta]$, where the identity is given by the “empty” polygonal loop 0 starting and ending at v , and with inverse given by $i(\gamma) = \gamma^{-1}$.

Exercise 3.10. Check that the above makes sense, and that $p_1(\Sigma, v)$ is a group.

Now suppose that Σ is obtained by glueing up the sides of a single polygon P in pairs. Suppose further that after the result of this glueing, all the vertices of this polygon are identified to a single vertex v . Let e_1, \dots, e_n be the edges and

$$\gamma = e_{i_1}^{\pm 1} e_{i_2}^{\pm 1} \dots e_{i_m}^{\pm 1}$$

the oriented boundary of P . Then there is a natural isomorphism

$$p_1(\Sigma, v) = \langle e_1, \dots, e_n | \gamma \rangle$$

that is, p_1 can be thought of as the group generated by the edges e_i subject to the relation defined by γ . Notice that each of the e_i appears twice in γ , possibly with distinct signs. If Σ is orientable, each e_i appears with opposite signs.

Example 3.11. Let T denote the surface obtained by glueing opposite sides of a square by translation. Thus the edges of the square can be labelled (in the circular ordering) by a, b, a^{-1}, b^{-1} and a presentation for the group is

$$p_1(T, v) = \langle a, b | [a, b] \rangle$$

It is not too hard to see that this is isomorphic to the group $\mathbb{Z} \oplus \mathbb{Z}$.

3.1.3. *Homotopy theory.* The following definition generalizes the notation of homotopy equivalence of maps $S^1 \rightarrow \Sigma$:

Definition 3.12. Two continuous maps $f_1, f_2 : X \rightarrow Y$ are *homotopic* if there is a map $F : X \times I \rightarrow Y$ satisfying $F(\cdot, 0) = f_1$ and $F(\cdot, 1) = f_2$. If $M \subset X$ and $N \subset Y$ with $f_1|_M = f_2|_M$ and $f_i(M) \subset N$, then f_1, f_2 are *homotopic relative to M* if a map F can be chosen as above with $F(m, \cdot)$ constant for every $m \in M$.

The set of homotopy classes of maps from X to Y is usually denoted $[X, Y]$. If these space have basepoints x, y the set of maps taking x to y modulo the equivalence relation of homotopy relative to x is denoted $[X, Y]_0$.

We can define a very important category whose objects are topological spaces and whose morphisms are *homotopy equivalence classes of continuous maps*. A refinement of this category is the category whose objects are topological spaces with basepoints, and whose morphisms are *homotopy equivalence classes of continuous maps relative to base points*.

Definition 3.13. The *fundamental group* $\pi_1(X, x)$ of an arbitrary topological space with a basepoint x as the group whose elements are homotopy classes of maps $(S^1, 0) \rightarrow (X, x)$, where multiplication is defined by $[c_1][c_2] = [c_1 * c_2]$. In the notation above, the elements of $\pi_1(X, x)$ correspond to elements of $[S^1, X]_0$.

Lemma 3.14. Let $f : X \rightarrow Y$ be a continuous map taking x to y . Then f induces a natural homomorphism $f_* : \pi_1(X, x) \rightarrow \pi_1(Y, y)$.

Definition 3.15. A path-connected space X is *simply-connected* if $\pi_1(X, x)$ is the trivial group.

Observe that a path-connected space is simply-connected if and only if every loop in X can be shrunk to a point.

3.1.4. *Simplicial approximation.* The following is known as the *simplicial approximation theorem*:

Theorem 3.16. Let K, L be two simplicial complexes. Then any continuous map $f : K \rightarrow L$ is homotopic to a simplicial map $f' : K' \rightarrow L$ where K' is obtained from K by subdividing simplices. Furthermore, if $C \subset K$ is a simplicial subset, and $f : C \rightarrow L$ is simplicial, then we can require f' to agree with f restricted to C .

Exercise 3.17. Using this theorem, show that by choosing $p = v$, every continuous map $c : S^1 \rightarrow \Sigma$ taking the base point to v is homotopic through basepoint-preserving maps to a simplicial loop $\gamma \subset \Sigma$ which begins and ends at v . Moreover, two such simplicial loops are homotopic if and only if they differ by a sequence of elementary moves. Thus there is a natural isomorphism $\pi_1(\Sigma, p) \cong p_1(\Sigma, v)$.

This is actually a very powerful observation: the group $\pi_1(\Sigma, p)$ is *a priori* very difficult to compute, but manifestly doesn't depend on a piecewise linear structure on Σ . On the other hand, $p_1(\Sigma, v)$ is easy to compute (or at least find a presentation for), but it is *a priori* hard to see that this group, up to isomorphism, doesn't depend on the piecewise linear structure.

Exercise 3.18. Suppose K is a simplicial complex. Let K^2 denote the union of the simplices of K of dimension at most 2. Use the simplicial approximation theorem to show that for any vertex v of K , $\pi_1(K, v) \cong \pi_1(K^2, v)$.

3.1.5. *Covering spaces.*

Definition 3.19. A space Y is a *covering space* for X if there is a map $f : Y \rightarrow X$ (called a *covering projection*) with the property that every point $x \in X$ has an open neighborhood U such that $f^{-1}(U)$ is a disjoint union of open sets $U_i \subset Y$, and f maps each U_i homeomorphically to U . The *universal cover* of a space X (if one exists) is a simply-connected space \tilde{X} which is a covering space of X .

An open neighborhood U of a point x of the kind provided in the definition is said to be *evenly covered* by its preimages $f^{-1}(U)$. If X is locally connected, we can assume that the open neighborhoods which are evenly covered are connected.

For a path $I \subset X$ we can find, for each point $p \in I$ an open neighborhood U_p of p which is evenly covered. Since I is *compact* (see appendix) only finitely many open neighborhoods are needed to cover I ; call these U_1, \dots, U_n . If we let V_1 denote some component of $f^{-1}(U_1)$ which maps homeomorphically to U_1 . Then there is a unique map $g : I \cap U_1 \rightarrow V_1$ such that $fg = \text{id}$. Moreover, there is a unique choice of V_2 from amongst the components of $f^{-1}(U_2)$ such that g can be extended to $g : I \cap (U_1 \cup U_2) \rightarrow V_1 \cup V_2$ with $fg = \text{id}$. Continuing inductively, we see that the choice of V_3, V_4, \dots, V_n are all uniquely determined by the original choice V_1 .

Exercise 3.20. *Modify the above argument to show that for every map $g : I \rightarrow X$ any every $p \in f^{-1}g(0)$ there is a unique lift $\tilde{g} : I \rightarrow Y$ such that $\tilde{g}(0) = p$ and $f\tilde{g} = g$.*

Exercise 3.21. *Suppose $g_1, g_2 : I \rightarrow X$ with $g_1(0) = g_2(0)$ and $g_1(1) = g_2(1)$ are homotopic through homotopies which keep the endpoints of I fixed. Let \tilde{g}_1 be some lift of g_1 . Show that the lift \tilde{g}_2 of g_2 with $\tilde{g}_1(0) = \tilde{g}_2(0)$ also satisfies $\tilde{g}_1(1) = \tilde{g}_2(1)$, and these two lifts are also homotopic rel. endpoints. (Hint: let $G : I \times I \rightarrow X$ be a homotopy between g_1 and g_2 . Try to “lift” G to a map $\tilde{G} : I \times I \rightarrow Y$ satisfying appropriate conditions.)*

Let $p \in X$ be a basepoint, and let $\tilde{p} \in f^{-1}(p)$. Then the projection $f : Y \rightarrow X$ induces a homomorphism $f_* : \pi_1(Y, \tilde{p}) \rightarrow \pi_1(X, p)$. Let $K \subset \pi_1(X, p)$ denote the image of f_* . Then a loop $\alpha : S^1 \rightarrow X$ with $[\alpha] \in K$ can be lifted to a loop $\tilde{\alpha} : S^1 \rightarrow Y$, by the argument of the previous exercise. Conversely, if two loops α, β represent the same element of K , then their lifts represent the same element of $\pi_1(Y, \tilde{p})$. Thus we may identify $\pi_1(Y, \tilde{p})$ with the subgroup K .

Exercise 3.22. *Show that for a space Z and a map $g : Z \rightarrow X$ there is a lift of g to $\tilde{g} : Z \rightarrow Y$ with $f\tilde{g} = g$ if and only if $g_*(\pi_1(Z)) \subset K$.*

Definition 3.23. A space is *semi-locally simply-connected* if every point p has a neighborhood U so that every loop in U can be shrunk to a point in a possibly larger open neighborhood V .

Informally, a space is semi-locally simply-connected if sufficiently small loops in the space are homotopically inessential.

Theorem 3.24. *Let X be connected, locally connected and semi-locally simply-connected. For any subgroup $G \subset \pi_1(X, x)$ there is a covering space X_G of X and a point $y \in f^{-1}(x)$ such that $f_*(\pi_1(X_G, y)) = G$.*

Sketch of proof: Let $\Omega(X)$ be the space of paths $\gamma : I \rightarrow X$ which start at x . We induce an equivalence relation on $\Omega(X)$ where we say two paths γ_1, γ_2 are equivalent if $[\gamma_2^{-1} * \gamma_1] \in G$. Set $X_G = \Omega(X) / \sim$. Since X is semi-locally simply connected, for

sufficiently small paths γ_1, γ_2 between points $x_1, x_2 \in X$ the loop $\gamma_2^{-1} * \gamma_1$ is contractible. So any two paths which differ only by substituting γ_1 in one for γ_2 in the other will be equivalent in $\Omega(X)$, and this says that the equivalence classes of $\Omega(X)$ are parameterized locally by points in X ; that is, X_G is a covered space of X . It can be verified that X_G is simply connected, and that it satisfies the conditions of the theorem. \square

Exercise 3.25. *Fill in the gaps in the sketch of the proof above.*

Exercise 3.26. *Show that the universal cover of a space X , if it exists, is unique, by using the lifting property.*

3.1.6. Discrete groups.

Definition 3.27. Let Γ be a group of symmetries of a space X which is one of $\mathbb{S}^n, \mathbb{E}^n, \mathbb{H}^n$ for some n . Γ is *properly discontinuous* if, for each closed and bounded subset K of X , the set of $\gamma \in \Gamma$ such that $\gamma(K) \cap K \neq \emptyset$ is finite. Γ acts *freely* if no $\gamma \in \Gamma$ has a fixed point; that is, if $\gamma(p) = p$ for any p , then $\gamma = \text{id}$.

For a point $p \in X$, the subgroup of Γ which fixes p is called the *stabilizer* of p , and is typically denoted by $\Gamma(p)$. If Γ acts properly discontinuously, then $\Gamma(p)$ is finite for any p .

Definition 3.28. A subgroup Γ of a Lie group G is *discrete* if $K \cap \Gamma$ is finite for any compact subset $K \subset G$.

Exercise 3.29. *Show that if G is a Lie group of symmetries of a space X as above, then a subgroup Γ is discrete iff it acts properly discontinuously on X .*

Suppose Γ acts on X freely and properly discontinuously. We can define a quotient space $M = X/\Gamma$ where two points $x, y \in X$ are identified exactly when there is some $\gamma \in \Gamma$ such that $\gamma(x) = y$.

Theorem 3.30. *If Γ acts on X freely and properly discontinuously, where X is one of $\mathbb{S}^n, \mathbb{E}^n, \mathbb{H}^n$ for some n , then the projection $X \rightarrow X/\Gamma$ is a covering space, and X is the universal cover of X/Γ .*

Proof: Pick a point $x \in X$ and let U be a neighborhood of x which intersects only finitely many translates $\gamma(U)$. Then there is a smaller $V \subset U$ which is disjoint from all its translates. Then under the quotient map, each translate $\gamma(V)$ is mapped homeomorphically to its image. Thus X is a covering space of X/Γ , and since it is simply-connected, it is the universal cover. \square

3.1.7. *Fundamental domains.* Let Γ act on X properly discontinuously, where X is one of the spaces $\mathbb{S}^n, \mathbb{E}^n, \mathbb{H}^n$.

Definition 3.31. A *fundamental domain* for the action of Γ is a polygon $P \subset X$ such that for all p in the interior of P , $\alpha(p) \cap P = \emptyset$ unless $\alpha = \text{id}$, and such that the faces of P are paired by the action of Γ .

The translates of a fundamental domain are disjoint except along their boundaries. Moreover, these translates cover all of X . Thus they give a tessellation of X , whose symmetry group contains Γ as a subgroup. For “generic” fundamental domains, there are no “accidental” symmetries, and the group of symmetries of the tessellation is exactly Γ . In general, fundamental domains can be *decorated* with some extra marking which destroys any additional extra symmetries of a fundamental domain.

Definition 3.32. Choose a point $p \in X$. The *Dirichlet domain* of Γ centered at p , is the set

$$D = \{q \in X \text{ such that } d(p, q) \leq d(p, \alpha(q)) \text{ for all } \alpha \in \Gamma\}$$

In dimension 2, 3, in the cases we are interested in, D will be a locally finite polygon in X whose faces are paired by the action of Γ ; but in general, D might not be a polygon.

A Dirichlet domain is a fundamental domain for Γ .

3.2. Lattices in \mathbb{E}^2 .

3.2.1. *Discrete groups in $\text{Isom}(\mathbb{E}^2)$.* Perhaps the most important theorem about discrete subgroups of $\text{Isom}(\mathbb{E}^n)$ is Bieberbach's theorem, that such a group has a free abelian subgroup of finite index. In dimension two, this can be refined as follows:

Theorem 3.33. *Let Γ act properly discontinuously on \mathbb{E}^2 by isometries. Then Γ has a subgroup Γ^+ of index at most 2 which is orientation-preserving. Moreover, Γ^+ is one of the following:*

- Γ^+ is a subgroup of $\text{stab}(p)$ for some p . In this case, $\Gamma^+ \cong \mathbb{Z}/n\mathbb{Z}$ consists of powers of a single rotation.
- Γ^+ is a semi-direct product

$$\Gamma^+ = \mathbb{Z}/n\mathbb{Z} \ltimes \mathbb{Z}$$

where the \mathbb{Z} factor is generated by a translation, and n is 1 or 2. In the second case, the conjugation action takes $x \rightarrow -x$.

- Γ^+ is a semi-direct product

$$\Gamma^+ = \mathbb{Z}/n\mathbb{Z} \ltimes (\mathbb{Z} \oplus \mathbb{Z})$$

where the $\mathbb{Z} \oplus \mathbb{Z}$ factor is generated by a pair of linearly independent translations, and $n = 1, 2, 3, 4$ or 6 . If $n = 4$, the $\mathbb{Z} \oplus \mathbb{Z}$ is conjugate to the group of translations of the form $z \rightarrow z + n + mi$ for integers n, m . If $n = 3$ or 6 , the $\mathbb{Z} \oplus \mathbb{Z}$ is conjugate to the group of translations of the form $z \rightarrow z + n + m \frac{1+i\sqrt{3}}{2}$ for integers n, m .

Proof: First, there is a homomorphism

$$o : \text{Isom}(\mathbb{E}^2) \rightarrow \mathbb{Z}/2\mathbb{Z}$$

where $o(\alpha)$ is 0 or 1 depending on whether or not α is orientation preserving. The intersection of Γ with the kernel of o is Γ^+ , and it has index at most 2. Next, there is a homomorphism

$$a : \text{Isom}^+(\mathbb{E}^2) \rightarrow S^1$$

given by the action of isometries on equivalence classes of parallel lines, where the image is the angles of rotation. There is an induced homomorphism

$$a : \Gamma^+ \rightarrow S^1$$

The kernel K of a in Γ^+ consists of a group of translations, and is therefore abelian.

Suppose $\theta = a(\alpha)$ for some $\alpha \in \Gamma^+$, and θ nonzero. Then α is a rotation, and therefore fixes some p . Since Γ^+ is properly discontinuous, either $\Gamma^+ \subset \text{stab}(p)$ in which case Γ^+ is cyclic and consists of the powers of a fixed rotation, or there is a (non-unique) closest image $q \neq p$ of p under some β . Since $q = \beta(p)$, no translate of p is closer to q than p . If $|\theta| < \frac{2\pi}{6}$, then

$$0 < d(q, \alpha(q)) = 2 \sin\left(\frac{\theta}{2}\right) d(p, q) < d(p, q)$$

which is a contradiction. Since the same must also be true for each power of α , the order of α is ≤ 6 . If the order is 5, then $0 < d(\alpha(q), \beta\alpha^{-1}\beta^{-1}(p)) < d(p, q)$, which is a contradiction. Hence the image $a(\Gamma^+)$ is $\mathbb{Z}/n\mathbb{Z}$ where $n = 1, 2, 3, 4$ or 6 . In any case, the image is cyclic, and is therefore generated by a single $a(\alpha)$ where α is a rotation, so Γ^+ is a semi-direct product.

The kernel K of a in Γ^+ is a properly discontinuous group of translations. If K is nontrivial and the elements of K are linearly dependent, they are all powers of the element α of shortest translation length. The image $a(\Gamma^+)$ must preserve the translation of α ; in particular, this image must be trivial or $\mathbb{Z}/2\mathbb{Z}$.

If K is nontrivial and the elements of K are linearly independent, they are generated by two elements of shortest translation length. The image $a(\Gamma^+)$ must preserve the set of nonzero elements of shortest length, so if this image is $\mathbb{Z}/4\mathbb{Z}$ the group K is generated by two perpendicular vectors of equal length. If the image is $\mathbb{Z}/3\mathbb{Z}$ or $\mathbb{Z}/6\mathbb{Z}$, K is generated by two vectors at angle $\frac{2\pi}{6}$ of equal length. \square

The two special lattices (i.e. groups of translations generated by two linearly independent elements) appearing in theorem 3.33 are often called the *square* and *hexagonal* lattices respectively.

3.2.2. *Integral quadratic forms.* A good reference for the material in this and the next section is [2].

Definition 3.34. A *quadratic form* is a homogeneous polynomial of degree 2 in some variables. That is, a function of variables x_1, \dots, x_n such that

$$f(\lambda_1 x_1, \dots, \lambda_n x_n) = \prod_{i=1}^n \lambda_i^2 f(x_1, \dots, x_n)$$

A quadratic form is *integral* if its coefficients as a polynomial are integers.

We will be concerned in the sequel with integral quadratic forms of two variables, such as $3x^2 + 2xy - 7y^2$ or $-x^2 - 3xy$.

Notice for every quadratic form $f(\cdot, \cdot)$ there corresponds uniquely a symmetric matrix M_f such that $f(x, y) = \begin{bmatrix} x & y \end{bmatrix} M_f \begin{bmatrix} x \\ y \end{bmatrix}$. In particular, if $f(x, y) = ax^2 + bxy + cy^2$ then

$$M_f = \begin{bmatrix} a & \frac{b}{2} \\ \frac{b}{2} & c \end{bmatrix}$$

Definition 3.35. We say that a quadratic form $f(x, y)$ *represents* an integer n if there is some assignment of integer values n_1, n_2 to x, y for which

$$f(n_1, n_2) = n$$

Given an integral quadratic form, it is a natural question to ask what interger values it represents. Observe that there are some elementary transformations on quadratic forms which do not change the set of values represented. If we substitute $x \rightarrow x \pm y$ or $y \rightarrow y \pm x$ we get a new quadratic form. Since this substitution is *invertible*, the new quadratic form obtained represents exactly the same set of values as the original. For instance,

$$x^2 + y^2 \longleftrightarrow x^2 + 2xy + 2y^2$$

This substitution defines an equivalence relation on quadratic forms.

The most general form of substitution allowed is a transformation of the form

$$x \rightarrow px + qy, \quad y \rightarrow rx + sy$$

For integers p, q, r, s . Again, since this substitution must be invertible, we should have $ps - qr = \pm 1$. That is, the matrix $\begin{bmatrix} p & q \\ r & s \end{bmatrix}$ is in $\pm SL(2, \mathbb{Z})$. Since these forms are homogeneous of degree 2, the substitution $x \rightarrow -x, y \rightarrow -y$ does nothing. One can easily check that every other substitution has a nontrivial effect on some quadratic form.

In particular, we have the following theorem:

Theorem 3.36. *Integral quadratic forms up to equivalence are parameterized by equivalence classes of matrices of the form $\begin{bmatrix} a & \frac{b}{2} \\ \frac{b}{2} & c \end{bmatrix}$ for integers a, b, c modulo the conjugation action of $PSL(2, \mathbb{Z})$.*

Observe that what is really going on here is that we are evaluating the quadratic form f on the integral lattice $\mathbb{Z} \oplus \mathbb{Z}$. The group $SL(2, \mathbb{Z})$ acts by automorphisms of this lattice, and therefore permutes the set of values attained by the form f .

There is nothing special about the integral lattice here; it is obvious that $SL(2, \mathbb{Z})$ acts by automorphisms of any lattice L . In particular, if $L = \langle e_1, e_2 \rangle$ then $M = \begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix} \in SL(2, \mathbb{Z})$ acts on elements of this lattice by

$$M \cdot re_1 + se_2 \rightarrow (\alpha r + \beta s)e_1 + (\gamma r + \delta s)e_2$$

3.2.3. Moduli of tori, continued fractions and $PSL(2, \mathbb{Z})$.

Definition 3.37. Let r be a real number. A *continued fraction expansion* of r is an expression of r as a limit of a (possibly terminating) sequence

$$n_1, n_1 + \frac{1}{m_1}, n_1 + \frac{1}{m_1 + \frac{1}{n_2}}, n_1 + \frac{1}{m_1 + \frac{1}{n_2 + \frac{1}{m_2}}}, \dots$$

where each of the n_i, m_i is a positive integer.

A continued fraction expansion of r can be obtained inductively by Euclid's algorithm. First, n_1 is the biggest integer $\leq r$. So $0 \leq r - n_1 < 1$. If $r - n_1 = 0$ we are done. Otherwise, $r' = \frac{1}{r - n_1} > 1$ and we can define m_1 as the biggest integer $\leq r'$. So $0 \leq r' - m_1 < 1$. continuing inductively, we produce a series of integers $n_1, m_1, n_2, m_2, \dots$ which is the *continued fraction expansion* of r . If r is rational, this procedure terminates at a finite stage. The usual notation for the continued fraction expansion of a real number r is

$$r = n_1 + \frac{1}{m_1 + \frac{1}{n_2 + \frac{1}{m_2 + \frac{1}{n_3 + \dots}}}}$$

The following theorem is quite easy to verify:

Theorem 3.38. *If n_1, m_1, \dots is a continued fraction expansion of r , then the successive approximations*

$$n_1, n_1 + \frac{1}{m_1}, n_1 + \frac{1}{m_1 + \frac{1}{n_2}}, \dots$$

denoted r_1, r_2, r_3, \dots satisfy

$$|r - r_i| \leq |r - p/q|$$

for any integers (p, q) , where $q <$ the denominator of r_{i+1} .

Thus, the continued fraction approximations of r are the best rational approximations to r for a given bound on the denominator.

Let T be a flat torus. Then the isometry group of T is transitive (this is not too hard to show). Pick a point p , and cut T up along the two shortest simple closed curves which start and end at p . This produces a Euclidean parallelogram P . After rescaling T , we can assume that the shortest side has length 1. We place P in \mathbb{E}^2 so that this short side is the segment from 0 to 1, and the other side runs from 0 to z where $\text{Im}(z) > 0$. By hypothesis, $|z| \geq 1$. Moreover, if $|\text{Re}(z)| \geq \frac{1}{2}$ we can replace z by $z + 1$ or $z - 1$ with smaller norm, contradicting the choice of curves used for the decomposition. Let D be the region in the upper half-plane bounded by the two vertical lines $\text{Re}(z) = \frac{1}{2}$, $\text{Re}(z) = -\frac{1}{2}$ and the circle $|z| = 1$.

The group $PSL(2, \mathbb{Z})$ acts naturally on \mathbb{H} as a subgroup of $PSL(2, \mathbb{R})$. The action there is properly discontinuous. The action of $PSL(2, \mathbb{Z})$ permutes the sides of D . The element $\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ pairs the two vertical sides, and the element $\begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$ preserves the bottom side, interchanging the left and right pieces of it. In particular, D is a fundamental domain for $PSL(2, \mathbb{Z})$. The quotient is topologically a disk, but with two ‘‘cone points’’ of order 2 and 3 respectively, which correspond to the points i and $\frac{1+i\sqrt{3}}{2}$ respectively, whose stabilizers are $\mathbb{Z}/2\mathbb{Z}$ and $\mathbb{Z}/3\mathbb{Z}$ respectively.

This quotient is an example of an *orbifold*.

Definition 3.39. An n -dimensional *orbifold* is a space which is locally modelled on \mathbb{R}^n modulo some finite group.

A 2-dimensional orbifold looks like a surface except at a collection of isolated points p_i where it looks like the quotient of a disk by the action of $\mathbb{Z}/n_i\mathbb{Z}$, a group of rotations centered at p_i . The point p_i is a *cone point*, sometimes also called an *orbifold point*. The finite group is part of the data of the orbifold. One can think of the orbifold combinatorially as a surface (in the usual sense) with a finite number of distinguished points, each of which has an integer attached to it. Geometrically, this point looks like a ‘‘cone’’ made from a wedge of angle $2\pi/n_i$.

We can define an *orbifold fundamental group* $\pi_1^o(\cdot)$ for a surface orbifold. Thinking of our orbifold Σ as X/Γ for the moment where Γ acts properly discontinuously but not freely, the orbifold fundamental group of Σ should be exactly Γ . This means that a small loop around an orbifold point p_i should have order n_i in $\pi_1^o(\Sigma)$. Note that we are being casual about basepoints here, so we are only thinking of these groups up to isomorphism.

In any case, the orbifold $\mathbb{H}^2/PSL(2, \mathbb{Z})$ should have orbifold fundamental group isomorphic to $PSL(2, \mathbb{Z})$. There is an element of order 2 corresponding to the loop around the order 2 point; a representative of this element in $PSL(2, \mathbb{Z})$ is $\begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$. There is an element of order 3 corresponding to the loop around the order 3 point; a representative of this element in $PSL(2, \mathbb{Z})$ is $\begin{bmatrix} 1 & 1 \\ -1 & 0 \end{bmatrix}$. Note that these elements have order 2 and 3 respectively in $PSL(2, \mathbb{Z})$, even though the corresponding matrices have orders 4 and 6 respectively in $SL(2, \mathbb{Z})$.

Every loop in the disk can be shrunk to a point; it follows that every loop in the orbifold $\mathbb{H}^2/PSL(2, \mathbb{Z})$ can be shrunk down to a collection of small loops around the two cone points in some order. That is, the group $PSL(2, \mathbb{Z})$ is generated by these two elements.

Theorem 3.40. *A presentation for $PSL(2, \mathbb{Z})$ is given by*

$$PSL(2, \mathbb{Z}) \cong \langle \alpha, \beta | \alpha^2, \beta^3 \rangle = \mathbb{Z}/2\mathbb{Z} * \mathbb{Z}/3\mathbb{Z}$$

where representatives of α and β are the two matrices given above.

Proof: By the discussion above, all that needs to be established is that there are no other relations that do not follow from the relations $\alpha^2 = \text{id}$ and $\beta^3 = \text{id}$. That is, there is a homomorphism $\phi : G \rightarrow PSL(2, \mathbb{Z})$ sending α, β to the two matrices given; all we need to check is that the kernel of this homomorphism consists of the identity element.

A general element of $G = \langle \alpha, \beta | \alpha^2, \beta^3 \rangle$ is a product $\alpha^{a_1} \beta^{b_1} \alpha^{a_2} \beta^{b_2} \dots \alpha^{a_n} \beta^{b_n}$ where each of the a_i, b_i are integers. We reduce the $a_i \pmod 2$ and the $b_i \pmod 3$; after rewriting of this kind, we are left with a product of the form above where every $a_i = 1$ and every b_i is 1 or 2. We write $L = \alpha\beta$ and $R = \alpha\beta^2$, so that every nontrivial element of G is of the form $w, \beta^{\pm 1}w, w\alpha, \beta^{\pm 1}w\alpha$ where w is a word in the letters L and R . Furthermore, we have the relation $(\alpha\beta\alpha\beta^2)^3$

We show that no word w in the letters L and R is trivial in the group $PSL(2, \mathbb{Z})$. A similar argument works for elements of G of the other forms.

Now, $\phi(L) = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$ and $\phi(R) = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ in $PSL(2, \mathbb{Z})$. Suppose that

$$w = L^{m_1} R^{n_1} L^{m_2} R^{n_2} \dots L^{m_k} R^{n_k}$$

where all the m_i, n_i are nonzero, say. Then we can calculate

$$\phi(w) = \begin{bmatrix} p & r \\ q & s \end{bmatrix}$$

where

$$\frac{p}{q} = \frac{1}{m_1 + n_1} \frac{1}{m_2 + n_2} \dots \frac{1}{m_k}$$

$$\frac{r}{s} = \frac{1}{m_1 + n_1} \frac{1}{m_2 + n_2} \dots \frac{1}{m_k + n_k}$$

where the notation is for a continued fraction expansion. That is, the alternating coefficients m_i, n_i give the continued fraction expansions of $\frac{p}{q}$ and $\frac{r}{s}$. In particular, $\phi(w) \neq \text{Id}$ unless w is the empty word, and ϕ is an isomorphism. \square

Notice actually that this method of proof does considerably more. We have shown that every element obtained by a product of *positive* multiples of L and R is non-trivial. It is not true that the *group* generated by L and R is free, since LR^{-1} has order 3. In fact, L and R together generate the entire group $PSL(2, \mathbb{Z})$. But the group generated by L^2 and R^2 is free, since a fundamental domain for its action is the domain D' bounded by the lines

$$\text{Re}(z) = 1, \text{Re}(z) = -1$$

and the semicircles

$$|z - 1/2| = 1/2, |z + 1/2| = 1/2$$

Let Γ denote this subgroup of $PSL(2, \mathbb{Z})$. The domain D' is a regular ideal hyperbolic quadrilateral; the quotient \mathbb{H}^2/Γ is therefore topologically a punctured torus. By the argument of the previous section, a presentation is

$$\pi_1(\text{punctured torus}) \cong \langle \alpha, \beta \rangle$$

That is, $\Gamma \cong \mathbb{Z} * \mathbb{Z}$.

Another description of Γ is the following: there is an obvious homomorphism $\psi : PSL(2, \mathbb{Z}) \rightarrow PSL(2, \mathbb{Z}/2\mathbb{Z})$ given by reducing the entries mod 2. The image group has order 6 and the surjection is onto, so the kernel is a subgroup of index 6. Since the fundamental domain D' can be made from 6 copies of D , it follows that the index of Γ in $PSL(2, \mathbb{Z})$ is 6. Moreover, Γ is certainly contained in the kernel of ψ . It follows that Γ is exactly equal to this kernel.

Γ is sometimes also denoted by $\Gamma(2)$ (for “reduction mod 2”) and is of considerable interest to number theorists, who like to refer to it as the *principal congruence subgroup of level 2*.

Notice that the domain D' is obtained from two ideal triangles. The union of all the translates of D' by $\Gamma(2)$ gives a tessellation of \mathbb{H}^2 by regular ideal quadrilaterals; a subdivision of D' into two ideal triangles gives a subdivision of \mathbb{H}^2 into ideal triangles. If we choose the subdivision along the line $\operatorname{Re}(z) = 0$, the ideal triangulation T of \mathbb{H}^2 so obtained admits reflection symmetry along every edge. The 1-skeleton of the dual cell-decomposition to this ideal triangulation is the *infinite 3-valent tree*. There is a natural action of $PSL(2, \mathbb{Z})$ on this tree, where the elements of order 3 are the stabilizers of vertices and the elements of order 2 are the stabilizers of edges. This description of $PSL(2, \mathbb{Z})$ as a group of automorphisms of a tree gives another way to see that it is isomorphic to $\mathbb{Z}/2\mathbb{Z} * \mathbb{Z}/3\mathbb{Z}$.

For a rational point p/q we can consider the straight line l perpendicular to the real axis given by $\operatorname{Re}(z) = p/q$. As this line l moves from ∞ to p/q it crosses through many different triangles of T , and therefore determines a word w in the letters R, L and their inverses. By induction, it is easy to show that the word w is of the form

$$w = L^{m_1} R^{n_1} L^{m_2} R^{n_2} \dots L^{m_k} R^{n_k}$$

where

$$\frac{p}{q} = \frac{1}{m_1 + n_1 + m_2 + n_2 + \dots + m_k}$$

An irrational point r determines an infinite word

$$w = L^{m_1} R^{n_1} L^{m_2} R^{n_2} \dots$$

where

$$r = \frac{1}{m_1 + n_1 + m_2 + n_2 + \dots}$$

is an infinite continued fraction expansion of r .

Notice that this word w is *eventually periodic* exactly when r is of the form $a + \sqrt{b}$ for rational numbers a, b .

3.3. Finite subgroups of $SO(3)$ and \mathbb{S}^3 .

3.3.1. *The “fair dice”.* A *die* is a convex 3-dimensional polyhedron. We can ask under what conditions a die is *fair* — that is, the probability that the die will land on a given side is $1/n$ where n is the number of sides. This is a very hard problem to treat in full generality, since it is very hard to calculate these probabilities for a generic polyhedron. But there are certain circumstances under which it is easy to show that these probabilities are all equal; if for any two faces f_1, f_2 of a die D there is a symmetry of D to itself taking f_1 to f_2 then the die is manifestly fair. The group G of all symmetries of D is a subgroup of the group of permutations of the vertices. Any symmetry of D extends to an isometry of \mathbb{E}^3 , in particular it is an affine map. It follows that if the vertices of D are at the vectors v_1, v_2, \dots, v_n then the images of these vertices under a symmetry σ are the same set of

vectors in permuted order. Thus the symmetry fixes the center of gravity of D ; as a vector this is $\frac{\sum_{i=1}^n v_i}{n}$.

Translating this center of gravity to the origin in \mathbb{E}^3 , we see that G is a finite subgroup of $O(3)$. That is, G is a properly discontinuous group of isometries of \mathbb{S}^2 .

3.3.2. *Spherical orbifolds.* Of course, any properly discontinuous group Γ of isometries of \mathbb{S}^2 has a subgroup Γ^+ of index at most two which consists of orientation-preserving elements. Every orientation-preserving isometry of \mathbb{S}^2 has a fixed point, so Γ^+ does not act freely unless it is trivial. In any case, the quotient \mathbb{S}^2/Γ^+ will be a *spherical orbifold* Σ . This orbifold is topologically a surface with finitely many cone points p_1, \dots, p_m of orders n_1, \dots, n_m . The Gauss–Bonnet formula gives

$$\text{area}(\Sigma) = \int_{\Sigma} \kappa = 2\pi \left(\chi(\Sigma) - \sum_{i=1}^m \frac{n_i - 1}{n_i} \right)$$

Since the area is positive, Σ must be topologically a sphere, since that is the only surface with positive Euler characteristic. Each $\frac{n_i - 1}{n_i}$ term is at least $1/2$, so it follows that there can be at most 3 cone points. Notice too that if Σ has two cone points, small loops around them are isotopic, and therefore should represent the same element of the orbifold fundamental group; in particular, they should have the same order. Similarly, Σ cannot have a single cone point, since a nontrivial element of the orbifold fundamental group could be shrunk to a point.

We therefore have the following theorem:

Theorem 3.41. *Let Γ be a properly discontinuous group of isometries of \mathbb{S}^2 . Then Γ has a subgroup Γ^+ of index at most 2 which is orientation-preserving. The following are the possibilities for Γ^+ :*

- Γ^+ fixes a pair of antipodal points. $\Gamma^+ \cong \mathbb{Z}/n\mathbb{Z}$ and is generated by a single rotation.
- Γ^+ is generated by two rotations r_1, r_2 of order 2 whose axes are at an angle of $\frac{2\pi}{n}$ to each other. The group $\Gamma^+ = \langle r_1, r_2 \rangle$ is the dihedral group D_n .
- The quotient \mathbb{S}^2/Γ^+ is a sphere with 3 cone points of orders $(2, 3, 3)$, $(2, 3, 4)$ or $(2, 3, 5)$. Γ^+ in these cases is the group of orientation-preserving symmetries of the regular tetrahedron, octahedron, and dodecahedron respectively. As a group, Γ^+ is isomorphic to A_4, S_4, A_5 respectively.

An orbifold Σ whose underlying surface is a sphere with three cone points p_1, p_2, p_3 is called a *triangle orbifold*. For the sake of generality, we can think of a puncture as a “cone point of order ∞ ”, so that $\mathbb{H}^2/PSL(2, \mathbb{Z})$ is the triangle orbifold with cone points of order $(2, 3, \infty)$. The triangle orbifold with cone points of order (p, q, r) will be denoted $\Delta(p, q, r)$.

A presentation for the triangle orbifold with cone points (p, q, r) is

$$\pi_1^o(\Delta(p, q, r)) \cong \langle \alpha, \beta \mid \alpha^p, \beta^q, (\alpha\beta)^r \rangle$$

Lemma 3.42. *For $r = 3, 4, 5$ there is an isomorphism*

$$\pi_1^o(\Delta(2, 3, r)) \rightarrow PSL(2, \mathbb{Z}/r\mathbb{Z})$$

where in each case, the image of the small loops α, β around the cone points of order 2, 3 correspond to the equivalence classes of matrices

$$\alpha \rightarrow \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, \beta \rightarrow \begin{bmatrix} 1 & 1 \\ -1 & 0 \end{bmatrix}$$

Proof: There are certainly homomorphisms from $\pi_1^o(\Delta(2, 3, r))$ onto $PSL(2, \mathbb{Z}/r\mathbb{Z})$ determined by the maps in question, since the relations $\alpha^2 = \text{id}$ and $\beta^3 = \text{id}$ hold in $PSL(2, \mathbb{Z}/n\mathbb{Z})$ for any n , and $\alpha\beta = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$ which has order r in $PSL(2, \mathbb{Z}/r\mathbb{Z})$.

To see that these maps are injective, observe that the orders of the groups for $r = 3, 4, 5$ are both equal to 12, 24 and 60 respectively, so the maps are isomorphisms. \square

For $r > 5$, the group $\pi_1^o(\Delta(2, 3, r))$ is infinite, and therefore cannot be isomorphic to $PSL(2, \mathbb{Z}/r\mathbb{Z})$. But for $r = \infty$ we have seen

$$\pi_1^o(\Delta(2, 3, \infty)) \cong PSL(2, \mathbb{Z})$$

The homomorphism $PSL(2, \mathbb{Z}) \rightarrow PSL(2, \mathbb{Z}/r\mathbb{Z})$ for $3 \leq r \leq 5$ is induced by the map from $\Delta(2, 3, \infty)$ to $\Delta(2, 3, r)$ which is the identity away from the special points, sends the order 2, 3 points to order 2, 3 points respectively, and ‘‘sends the puncture’’ to the cone point of order r .

3.3.3. Reflection groups, Coxeter diagrams. If P is a polyhedron in X one of $\mathbb{S}^n, \mathbb{E}^n, \mathbb{H}^n$ whose dihedral angles between top dimensional faces are all of the form π/m_i for integers m_i , the group G_P generated by reflections in these faces of P acts properly discontinuously on X with fundamental domain P . G_P has a subgroup of index 2 consisting of orientation preserving elements, which has as fundamental domain a copy of P and its mirror image P' . This follows from a theorem called *Poincaré’s polyhedron theorem*. A precise statement and discussion are found in [7].

If $X = \mathbb{S}^n$, we can think of $\mathbb{S}^n \subset \mathbb{R}^{n+1}$ as the unit sphere, and reflections through hyperplanes in \mathbb{S}^n correspond to reflections in hyperplanes through the origin in \mathbb{R}^{n+1} . If π_i, π_j are two of these hyperplanes, and the corresponding reflections are denoted r_i, r_j then the composition r_i, r_j is a rotation through an angle $2\theta_{ij}$, where θ_{ij} is the angle between the planes π_i and π_j , and the rotation is in the plane spanned by the two perpendiculars to π_i, π_j . A presentation for the group G generated by reflections in the sides of P is

$$G_P \cong \langle r_1, r_2, \dots, r_n | (r_1)^2, (r_2)^2, \dots, (r_n)^2, (r_1 r_2)^{m_{12}}, \dots, (r_i r_j)^{m_{ij}}, \dots \rangle$$

where the angle between π_i and π_j is π/m_{ij} .

The subgroup G_P^+ of orientation-preserving elements of G_P is generated by the elements of the form $r_i r_j$, which is to say, G_P^+ consists of products of even numbers of reflections.

Definition 3.43. A *Coxeter group* G is an abstract group defined by a group presentation of the form $\langle r_i | (r_i r_j)^{k_{ij}} \rangle$ where

- the indices vary over some index set I
- the exponent $k_{ij} = k_{ji}$ is either a positive integer or ∞ for each pair i, j
- $k_{ii} = 1$ for each i
- $k_{ij} > 1$ for each $i \neq j$

If $k_{ij} = \infty$ for some i, j then the corresponding relation is meaningless and may be deleted from the presentation.

Definition 3.44. The *Coxeter graph* of the Coxeter group G is a labelled graph Γ with vertices corresponding to the index set I and edges $\langle (i, j) : k_{ij} > 2 \rangle$ labelled by k_{ij} .

For simplicity, edges with $k_{ij} = 3$ are usually left unlabelled.

Theorem 3.45. *Finite Coxeter groups can be realized as properly discontinuous spherical reflection groups.*

Notice that fundamental groups of triangle orbifolds are index 2 subgroups of reflection groups whose Coxeter graphs have three vertices.

3.3.4. “Bad” orbifolds. If Σ is a spherical orbifold with two cone points of order $p, q > 1$ where $p \neq q$, the orbifold Euler characteristic of Σ is $2 - \frac{p-1}{p} - \frac{q-1}{q} > 0$, so the universal cover of Σ should be \mathbb{S}^2 . But we have seen that this is impossible; the loop around the point of order p is freely homotopic to the loop of order q , so a presentation for $\pi_1^o(\Sigma)$ is $\langle \alpha | \alpha^p, \alpha^q, \alpha^{p-q} \rangle$. This group is $\mathbb{Z}/d\mathbb{Z}$ where d is the greatest common divisor of p and q ; but every $\mathbb{Z}/d\mathbb{Z}$ subgroup of $SO(3, \mathbb{R})$ has as quotient the spherical orbifold with two cone points of order d .

Thus Σ is not obtained from a smooth surface by the action of a properly discontinuous group. An orbifold with no manifold cover is called a *bad orbifold*. A spherical orbifold with two cone points of unequal order is a bad orbifold; similarly, a spherical orbifold with one cone point is a bad orbifold.

It turns out that any orbifold whose underlying surface is not the sphere, but a surface of higher genus, is a good orbifold, and is obtained as a quotient X/Γ for X one of $\mathbb{S}^2, \mathbb{E}^2, \mathbb{H}^2$ and Γ a properly discontinuous group of isometries.

3.4. Discrete subgroups of $PSL(2, \mathbb{R})$.

3.4.1. *Glueing hyperbolic polygons.* Glueing up hyperbolic polygons to make a closed hyperbolic surface is not essentially different from glueing up spherical or flat polygons. If Σ_g denotes the unique (up to homeomorphism) closed orientable surface of genus $g > 1$, then we can decompose Σ_g (nonuniquely) into *pairs of pants*.

Definition 3.46. A *pair of pants* is the topological surface obtained from a disk by removing two subdisks — that is, a disk with two holes.

A pair of pants can also be thought of as a sphere minus three subdisks. The Euler characteristic of a pair of pants is -1 . Since the Euler characteristic of its boundary is 0, a surface obtained from glueing n pairs of pants has Euler characteristic $-n$. So Σ_g can be decomposed (in many different ways) into $2g - 2$ pairs of pants.

Exercise 3.47. *Show that the number of decompositions of Σ_g into pairs of pants, up to combinatorial equivalence, is equal to the number of graphs with $2g - 2$ vertices with 3 edges at every vertex. Such graphs are called trivalent graphs. Show that the number of such graphs is positive for $g > 1$, and enumerate such graphs for $g \leq 5$ (you might need to write a computer program . . .)*

For α a closed loop in Σ_g , a choice of hyperbolic metric on Σ_g determines a unique shortest loop α_g — a geodesic — which is homotopic to α . For, if $p \in \alpha$ and $\tilde{\alpha}$ denotes an arc in \mathbb{H}^2 , the universal cover of Σ_g , whose endpoints project to p and such that $\tilde{\alpha}$ projects to α under the covering map, then there is a unique isometry $\gamma \in PSL(2, \mathbb{R})$ corresponding to an element of $\pi_1(\Sigma_g)$ taking one end of $\tilde{\alpha}$ to the other. If $\tilde{\alpha}_g$ denotes the invariant axis of γ , then $\tilde{\alpha}_g/\gamma = \alpha_g$ a geodesic in Σ_g .

Exercise 3.48. *If α is an essential simple closed curve in Σ_g — that is, it is embedded and does not bound a disk, then α_g is also simple. Furthermore, if α, β are disjoint essential simple closed curves, their geodesic representatives α_g, β_g are disjoint.*

By the exercise, a combinatorial decomposition of Σ_g into pairs of pants determines, for a hyperbolic metric on Σ_g , a (combinatorially equivalent) decomposition of that surface into hyperbolic pairs of pants with geodesic boundary. Call such an object a *geodesic pair of pants*.

Let P be a geodesic pair of pants with boundary circles α, β, γ , and δ an embedded arc joining two distinct boundary components α, β . Then we can let Q be the surface, topologically a torus with two disks removed, obtained from two copies of P with opposite orientations glued along the pairs of circles corresponding to α, β . Then the two copies of δ make up a closed loop $\hat{\delta}$ which has a unique geodesic representative $\hat{\delta}_g \subset Q$. There is an orientation-reversing map i from Q to itself which fixes $\alpha \cup \beta$. By uniqueness, $\hat{\delta}_g$ is invariant under i , and therefore it intersects the boundary curves in right angles. We obtain an arc δ_g in P perpendicular to α and β . There are two other arcs ϵ_g, λ_g in P perpendicular to α, γ and β, γ . These decompose P into two right angled hyperbolic hexagons H_1, H_2 . The alternate sides of H_1 and H_2 are equal, and therefore they are isometric, by an orientation-reversing isometry.

Exercise 3.49. *Prove the claim made in the previous paragraph. That is, show that a right-angled hyperbolic hexagon is determined up to isometry by the lengths of three nonadjacent sides. Conversely show that for any three number $p, q, r > 0$ there is a right-angled hexagon with three nonadjacent sides with those lengths.*

In short we have proved the following fact:

Lemma 3.50. *Let Σ_g be a surface of genus g . A combinatorial decomposition into pairs of pants and a hyperbolic metric on Σ_g determine a decomposition of Σ_g into $2g - 2$ geodesic pairs of pants. The geometry of these pairs of pants is determined uniquely by the lengths of the closed geodesics along which Σ_g was decomposed.*

It remains to understand how the pairs of pants can be put back together to give Σ_g . For P_1, P_2 a pair of geodesic pairs of pants with boundary components $\alpha_1 \subset \partial P_1$ and $\alpha_2 \subset \partial P_2$ with the same length, for any two points $p_1 \in \alpha_1$ and $p_2 \in \alpha_2$ there is a unique way to glue P_1 to P_2 by identifying α_1, α_2 so that p_1, p_2 match up. There is a 1-parameter family of glueings, parameterized by the amount of “twisting” of these geodesics. In particular, the geometry of Σ_g is determined uniquely by the $3g - 3$ lengths of the geodesics along which it is decomposed, together with $3g - 3$ twist parameters.

Thus we have a correspondence:

$$(\text{metric on } \Sigma_g, \text{ pair of pants decomposition}) \longleftrightarrow (\mathbb{R}^+)^{3g-3} \times (\mathbb{R}/\mathbb{Z})^{3g-3}$$

Here the pair of pants decomposition is thought of as *ordered*, in the sense that the $3g - 3$ curves are given specific labels, which correspond to the $3g - 3$ co-ordinates on the right.

Although this is a nice characterization, the information contained in a pair of pants decomposition is both too little and too much — too little because we have not resolved the \mathbb{Z} -ambiguity in the twist parameters, and too much because it does not answer the question of what the space of hyperbolic structures on a surface is parameterized by. We address these issues now.

Definition 3.51. Fix a base surface Σ of genus $g > 1$. The space of *marked hyperbolic structures on Σ* , denoted $\mathcal{MH}(\Sigma)$ is the space of equivalence classes of pairs (f, Σ') where Σ' is a hyperbolic surface and $f : \Sigma \rightarrow \Sigma'$ is a homeomorphism, and two such pairs (f_1, Σ_1) and (f_2, Σ_2) are equivalent if there is an isometry $i : \Sigma_1 \rightarrow \Sigma_2$ such that $f_1 \circ i$ is homotopic to f_2 as a map from Σ to Σ_2 .

Exercise 3.52. Show that the relation defined in the definition of marked hyperbolic structure is really an equivalence relation. That is, show it is symmetric, reflexive and transitive.

Theorem 3.53. For Σ a surface of genus g , there is a 1–1 correspondence

$$\mathcal{MH}(\Sigma) \longleftrightarrow (\mathbb{R}^+)^{3g-3} \times \mathbb{R}^{3g-3}$$

The correspondence is defined as follows: there is a pair of pants decomposition along essential simple closed curves $\alpha_1, \dots, \alpha_{3g-3}$ for Σ , and a collection of loops $\beta_1, \dots, \beta_{3g-3}$ transverse to the α_i such that if (f, Σ') is an element in $\mathcal{MH}(\Sigma)$, the corresponding coordinates are given by

$$(\text{length}((f(\alpha_1))_g), \dots, \text{length}((f(\alpha_{3g-3}))_g), \text{twist}((f(\alpha_1))_g), \dots, \text{twist}((f(\alpha_{3g-3}))_g))$$

where the twist parameters are normalized so that twist 0 corresponds, for fixed lengths of $(f(\alpha_i))_g$, to the unique marked surface for which the length of β_i is minimized.

This requires some explanation. The image of a fixed pair of pants decomposition of Σ under f determines one in Σ' , and therefore the lengths of the decomposed geodesics are well-defined and the twist parameters are well-defined mod 2π . Let P_1, P_2 be two geometric pairs of pants glued along boundaries to make a sphere with 4 disks removed Q . If α_i is the common loop in $\partial P_1 \cap \partial P_2$, then β_i is a dual curve which cuts Q into two other pairs of pants P'_1, P'_2 such that P'_1 has one boundary component in common with each of P_1, P_2 and similarly for P'_2 . Then twisting α_i through 2π replaces β_i with a new curve $t_{\alpha_i}(\beta_i)$, the curve obtained by a *Dehn twist around* α_i . Briefly: β_i decomposes into two arcs δ, ϵ along α_i , and α_i decomposes into two arcs μ, ν along β_i . Then $t_{\alpha_i}(\beta_i)$ is the simple closed curve homotopic to $\delta * \mu * \epsilon * \nu$, where $*$ denotes concatenation of arcs. Imagine β_i as a rubber band on the surface Q . When the two sides P_1, P_2 are twisted independently along α_i , the rubber band becomes twisted up, and when P_1 and P_2 return to their original configuration, the rubber band detects how many full rotations the two sides went through. It is true, though we don't prove it here, that there is a unique rotation for which the length of the geodesic representative of β_i is minimized. For a reference, see [1]. Thus for fixed sets of lengths of the geodesic representatives of all the α_j , we have well-defined twist parameters which detect the amount of twisting relative to this minimal twist. This shows the map to parameter space is well-defined. Conversely, such a collection of parameters defines a collection of geodesic pairs of pants and instructions for glueing them together to give a marked hyperbolic structure on Σ . Thus the two sets are the same and the theorem is proved.

Definition 3.54. Let Σ be a closed surface. The *mapping class group* of Σ , denoted $\text{MC}(\Sigma)$, is defined to be the quotient group

$$\text{MC}(\Sigma) = \text{Homeo}(\Sigma) / \text{Homeo}_0(\Sigma)$$

where $\text{Homeo}_0(\Sigma)$ denotes the normal subgroup of self-homeomorphisms of Σ which are homotopic to the identity.

Exercise 3.55. Show $\text{Homeo}_0(\Sigma)$ is a normal subgroup of $\text{Homeo}(\Sigma)$.

Notice that $\text{Homeo}(\Sigma)$ acts on $\mathcal{MH}(\Sigma)$ by

$$\psi : (f, \Sigma') \rightarrow (\psi \circ f, \Sigma')$$

Moreover, if $\psi \in \text{Homeo}_0(\Sigma)$, then

$$\psi(f, \Sigma') \sim (f, \Sigma')$$

with respect to the equivalence relation defined on representatives. That is, $\text{MC}(\Sigma)$ acts on $\mathcal{MH}(\Sigma)$. Moreover, the quotient space is exactly the space of equivalence classes of elements in $\mathcal{MH}(\Sigma)$ where $(f_1, \Sigma_1) \sim (f_2, \Sigma_2)$ if and only if there is an isometry $i : \Sigma_1 \rightarrow \Sigma_2$. That is, two marked hyperbolic structures have the same orbits under $\text{MC}(\Sigma)$ if and only if the underlying hyperbolic structures (forgetting the marking) are equivalent. In particular, there is a corresponding action of $\text{MC}(\Sigma)$ on $(\mathbb{R}^+)^{3g-3} \times \mathbb{R}^{3g-3}$ and therefore a correspondence

$$\text{hyperbolic structures on } \Sigma \longleftrightarrow \{(\mathbb{R}^+)^{3g-3} \times \mathbb{R}^{3g-3}\} / \text{MC}(\Sigma)$$

The action of $\text{MC}(\Sigma)$ on \mathbb{R}^{6g-6} is properly discontinuous, but it is not free. Thus the space of hyperbolic structures on Σ is best thought of as an orbifold. This quotient space is also known as *moduli space*.

Exercise 3.56. *Verify the claims made above. In particular, show that the action of $\text{MC}(\Sigma)$ on $\mathcal{MH}(\Sigma)$ is well-defined, independently of the choice of representative of an element in $\mathcal{MH}(\Sigma)$.*

The space $\mathcal{MH}(\Sigma)$ is also known as the *Teichmüller space* of Σ , and denoted $\text{Teich}(\Sigma)$.

Definition 3.57. Let Σ be a closed surface.

Let $\text{Homeo}^+(\Sigma)$ denote the subgroup of Σ consisting of *orientation-preserving* homeomorphisms. Then define

$$\text{MC}^+(\Sigma) = \text{Homeo}^+(\Sigma) / \text{Homeo}_0(\Sigma)$$

Notice that in this definition we use implicitly the fact that for a *closed* surface, the subgroup of self-homeomorphisms homotopic to the identity are all orientation-preserving. This is not true for surfaces with boundary without some extra constraints on the boundary behaviour of these homeomorphisms.

Exercise 3.58. *Let Σ be the unit disk. Find a self-homeomorphism homotopic to the identity which is orientation-reversing. Do the same with Σ an annulus. What about if Σ is a punctured surface of genus $g \geq 2$?*

3.5. Dehn twists and Lickorish's theorem.

Definition 3.59. An oriented (polyhedral) simple closed curve c in a surface Σ and an annulus neighborhood A of c parameterized as $S^1 \times I$ define a homeomorphism $t_c : \Sigma \rightarrow \Sigma$ by

$$t_c : \begin{cases} x \rightarrow x & \text{for } x \text{ outside } A \\ (\theta, t) \rightarrow (\theta - 2t\pi, t) & \text{for } (\theta, t) \in A \end{cases}$$

This homeomorphism is known as a *Dehn twist* about c . As an element of $\text{MC}(\Sigma)$, it depends only on the isotopy class of c .

Note that $[t_c]^{-1} = [t_{c'}]$ where c' denotes c with the opposite orientation.

Exercise 3.60. *If $h : \Sigma \rightarrow \Sigma$ is a homeomorphism, p a simple closed loop in Σ , and $h(p) = q$, then*

$$t_p = h^{-1}t_qh$$

The following theorem is proved in [4], and is often referred to as the *Lickorish twist theorem*:

Theorem 3.61 (Lickorish). *If Σ_g denotes the oriented surface of genus g , then the group $\text{MC}^+(\Sigma_g)$ is generated by Dehn twists in $3g - 1$ (explicitly described) simple closed curves. In particular, this group is finitely generated.*

Sketch of proof: The method of proof proceeds as follows: let c be a simple closed curve in Σ , and let A be a collection of simple closed curves in Σ . Then either c intersects each element α of A not at all, exactly once, or exactly twice with opposite orientations, or there exists a loop d which intersects α fewer times than c , and each element of A at most as many times as c , such that $t_\alpha(c)$ has fewer intersections with α and the same number or fewer intersections with each other element of A . Proceeding inductively, we see that if C is a maximal collection of disjoint essential simple closed curves and ψ is a homeomorphism of Σ , then there are a sequence of Dehn twists t_1, t_2, \dots such that $t_n t_{n-1} \dots t_1 \psi(C)$ intersects C in one of finitely many possibilities. After twisting some more in elements of C , we can assume the image of C is one of finitely many possibilities, which can be explicitly identified. In short, ψ can be written as a product of Dehn twists.

Now, for each such twist t_c , we can replace t_c by $t_d t_{t_d(c)} t_d^{-1}$ where each $d, t_d(c)$ intersect C more simply than c . In this way, each t_c can be expanded as a product of Dehn twists in curves which intersect C very simply. After twisting in C , it follows that these involve only finitely many possibilities, which can be explicitly enumerated. \square

Remark 3.62. Casson has shown that the number of twist generators required is at most $2g + 2$. Furthermore, it is known that $\text{MC}^+(\Sigma)$ is generated by only 2 elements (which are not Dehn twists).

4. APPENDIX — WHAT IS GEOMETRY?

Geometry is a beast that can be approached from many angles. Four of the most important concepts that arise from our different primitive intuitions of geometry are *symmetry, measurement, analysis, and continuity*. We briefly discuss these four faces of geometry, and mention some fundamental concepts in each. Don't worry if these concepts seem very technical or abstract — think of this section as an abstraction of the concrete notions found in the main body of the text.

4.1. Klein's "Erlanger Programm". At an address at Friedrich–Alexander–Universitaet in Erlangen Germany on December 17 1872, Felix Klein proposed a program to unify the study of geometry by the use of algebraic methods, more specifically, by the use of *group theory*. In particular, the geometrical properties of a space can be understood and explored by a study of the *symmetries* of that space. These symmetries can be organized into a natural algebraic object — a group. Conversely, this group can often be given a natural geometric structure, and investigated in its turn as a geometric space! The interplay between geometry and algebra leads to an enrichment of both structures.

4.1.1. Category theory.

Example 4.1. This is not really an example, but rather a *template* for the examples we will meet that fit into Klein's program. We are given a space X together with some sort of structure. A *structure-preserving map* from X to itself is called a *morphism*. The map from X to itself which does nothing is a distinguished morphism, the *identity morphism*, denoted $\mathbf{1}_X$. A morphism f is *invertible* if there is another morphism f^{-1} such that $f \circ f^{-1} = f^{-1} \circ f = \mathbf{1}_X$. The invertible morphisms are also called *automorphisms*. The set of automorphisms of X is a group called $\text{Aut}(X)$, with $\mathbf{1}_X$ as the identity, and composition as multiplication. Observe that a structure on a space can be *defined* by the admissible morphisms. This is a simple example of what is known as a *category*; in particular, it is a category with one object X .

Definition 4.2. More generally, a *category* can be thought of as a collection of *objects* (denoted \mathcal{O}) and a collection of *morphisms* or admissible maps between objects (denoted \mathcal{M}). Every morphism m has a *source* object $s(m)$ and a *target* object $t(m)$, which might be the same object. For every object o , there is a special morphism called the *identity* morphism

$$\mathbf{1}_o : o \rightarrow o$$

which acts like the usual identity: i.e.

$$\mathbf{1}_x m = m \text{ for any } m \text{ with } t(m) = x$$

$$m \mathbf{1}_x = m \text{ for any } m \text{ with } s(m) = x$$

The composition of two morphisms is another morphism, and this composition is associative; composition can be expressed as a function $c : \mathcal{M} \times \mathcal{M} \rightarrow \mathcal{M}$. That is, c satisfies

$$c(m, c(n, r)) = c(c(m, n), r) \text{ for any } m, n, r \in \mathcal{M}$$

Sometimes a category is written as a 5-tuple $(\mathcal{O}, \mathcal{M}, s, t, c)$, but in practice it is frequently sufficient to specify the objects and the morphisms.

A category is something like a *class* in an object-oriented programming language like C++; one defines at the same time the *data types* (the objects in the category) and the *admissible functions* which operate on them (the morphisms).

Example 4.3. The category whose objects are all sets and whose morphisms are all functions between sets is a category called **SET**. If X is an object in **SET** (i.e. a set) then $\text{Aut}(X)$ is the group of permutations of X .

Example 4.4. The category whose objects are all groups and whose morphisms are all homomorphisms between groups is called **GROUP**.

A very readable introduction to category theory, with numerous exercises, are the notes by John Stallings [9].

4.2. Metric geometry. One of our basic intuitions in geometry is that of *distance*. In fact the word geometry literally means “measuring the earth”. Metric geometry is the study of the concept of distance, and its various generalizations and abstractions. A beautiful (but quite advanced) reference for this subject is [3].

4.2.1. *Metric spaces.*

Definition 4.5. A *metric space* X, d is a set X together with a function

$$d : X \times X \rightarrow \mathbb{R}_0^+$$

where \mathbb{R}_0^+ denotes the non-negative real numbers, with the following properties:

(1) d is *symmetric*. That is,

$$d(x, y) = d(y, x)$$

(2) d is *nondegenerate* in the sense that

$$d(x, y) = 0 \text{ iff } x = y$$

(3) d satisfies the *triangle inequality*. That is,

$$d(x, y) + d(y, z) \geq d(x, z)$$

for all triples $x, y, z \in X$.

Example 4.6. The real line \mathbb{R} is a metric space with

$$d(x, y) = |x - y|$$

Example 4.7. The plane \mathbb{R}^2 is a metric space with

$$d((x_1, y_1), (x_2, y_2)) = (x_1 - x_2)^2 + (y_1 - y_2)^2$$

Example 4.8. The plane \mathbb{R}^2 is a metric space with

$$d((x_1, y_1), (x_2, y_2)) = |x_1 - x_2| + |y_1 - y_2|$$

This metric is known as the *Manhattan metric*. Can you see why?

Definition 4.9. An *isometry* of a metric space X is a 1–1 and onto transformation of X to itself which preserves distances between points. The set of *isometries* of a space X is a group $\text{Isom}(X)$, where multiplication in the group is composition of symmetries, and e is the trivial symmetry which fixes every x in X . This is an example of a group of the form $\text{Aut}(X)$ where the relevant structure on X is that of the *category of metric spaces MET* whose objects are metric spaces and whose morphisms are isometries.

Definition 4.10. Isometries are frequently too restrictive for many circumstances; a typical metric space of study might admit no non–trivial isometries at all. We can enrich the structure by allowing as morphisms those maps which, though they don’t literally *preserve* distances between points, at least don’t increase distances between points by too much. Such a map is called a *Lipschitz map*, and metric spaces with these as morphisms define a category **LIP** which is in many ways a much more interesting object than **MET**.

A map $f : X \rightarrow Y$ between metric spaces is *bilipschitz* if there is a $K > 1$ so that

$$\frac{1}{K}d_Y(f(x), f(y)) \leq d_X(x, y) \leq Kd_Y(f(x), f(y))$$

One may think of this as a map which only distorts distances up to a bounded factor. A bilipschitz map is 1–1, since metrics are nondegenerate. An invertible Lipschitz map with Lipschitz inverse is bilipschitz, so that the automorphisms in the category **LIP** are bilipschitz. Moreover, the composition of two bilipschitz maps is bilipschitz.

- Exercise 4.11.**
- (1) Show that the set of bilipschitz self–maps is a group for $X = \mathbb{R}$ with the Euclidean metric.
 - (2) (Harder) Show that the set of bilipschitz self–maps is a group for $X = \mathbb{R}^2$ with the Euclidean metric, and also with the Manhattan metric.
 - (3) Show that the bilipschitz self–maps of \mathbb{R}^2 with the Euclidean or the Manhattan metric are the same

4.3. Differential geometry. Differential geometry is the abstraction of calculus and analysis on n –dimensional Euclidean space to generalized geometric spaces called “smooth Riemannian manifolds”. These are spaces which look like Euclidean space on a small scale, but on larger scales they are deformed or “curved”. Einstein’s theory of general relativity says that our own universe is a certain kind of curved space, where the curvature is proportional to the strength of the gravitational force; on the human scale it looks Euclidean, but near massive objects like neutron stars, the “curvature” of the space is evidenced by the bending of light rays. The concept of curvature is a very important connection between geometry and topology.

Since calculus and analysis are basically *local*, one can do calculus on such spaces, since on smaller and smaller scales they look more and more like \mathbb{R}^n so that limits, derivatives, differentiability etc. all make sense, and the tools of multivariable calculus can be transplanted to this setting.

The morphisms which preserve the structure used to do differential geometry are the *smooth maps*, generalizations of differentiable functions.

4.3.1. *Smooth Manifolds.* A *manifold* is a space which, on a small scale, resembles Euclidean space of some dimension. The dimension is usually assumed to be constant over the space, and is called the *dimension of the manifold*.

A circle or a line is an example of a 1–dimensional manifold. A sphere or the surface of a donut is an example of a 2–dimensional manifold. Our universe, or the space outside a knot or link are examples of 3–dimensional manifolds.

Definition 4.12. A *smooth manifold* is a manifold on which one can do Calculus. One covers the space with a collection of little snapshots called “charts” which are meant to be all the different possible choices of local parameters for the space. Technically one has a collection of charts, which are subsets U_i of the manifold M , and a collection of ways of parameterizing these charts as subsets of Euclidean space; that is, maps $\phi_i : U_i \rightarrow V_i$ which are continuous and have continuous inverses, where V_i is some open region in some \mathbb{R}^n . These maps should be compatible, in the sense that if two charts U_i, U_j overlap, the map $\rho = \phi_j \phi_i^{-1}$ between the appropriate subsets of \mathbb{R}^n (where ρ is defined) should be smooth (i.e. it should have continuous partial derivatives of all orders), and it should be *locally invertible*; that is, the matrix of partial derivatives

$$d\rho = \begin{bmatrix} \frac{\partial \rho(x_1)}{\partial x_1} & \frac{\partial \rho(x_2)}{\partial x_1} & \cdots & \frac{\partial \rho(x_n)}{\partial x_1} \\ \frac{\partial \rho(x_1)}{\partial x_2} & \frac{\partial \rho(x_2)}{\partial x_2} & \cdots & \frac{\partial \rho(x_n)}{\partial x_2} \\ \cdots & \cdots & \cdots & \cdots \\ \frac{\partial \rho(x_1)}{\partial x_n} & \frac{\partial \rho(x_2)}{\partial x_n} & \cdots & \frac{\partial \rho(x_n)}{\partial x_n} \end{bmatrix}$$

should be invertible at every point.

This definition seems bulky but it is actually quite elegant. When doing multivariable calculus, we are used to switching back and forth between local co–ordinates which might only be defined on certain subsets of \mathbb{R}^n .

Example 4.13. In the plane \mathbb{R}^2 , we might switch between x, y Cartesian co–ordinates and r, θ polar co–ordinates. Note that θ is not really a “co–ordinate” on the whole of \mathbb{R}^2 , since its value is only well–defined up to multiples of 2π , and at the origin there is no sensible value for it. These co–ordinates are actually maps from subsets of the manifold \mathbb{R}^2 to the “standard” Euclidean space, which in this case also happens to be \mathbb{R}^2 . The first chart U_1 can be taken to be all of \mathbb{R}^2 , and the function ϕ_1 is just the identity $\phi_1 : (x, y) \rightarrow (x, y)$. The second function ϕ_2 is not defined on all of \mathbb{R}^2 , and might be given in the chart $U_2 = \{x, y > 1\}$ for instance, by

$$\phi_2 : (x, y) \rightarrow \left(\sqrt{x^2 + y^2}, \arctan \frac{x}{y} \right)$$

Defining $\rho = \phi_2 \phi_1^{-1} = \phi_2$ as above, check that $d\rho$ is invertible everywhere in the overlap of the two charts.

The charts on a smooth manifold are just the collection of all the possible local co–ordinates for subsets of the manifold; this collection of charts is called an *atlas*.

Example 4.14. The spaces \mathbb{R}^n are smooth manifolds for any n .

Example 4.15. An open subset of a smooth manifold is itself a smooth manifold, by restriction of charts and functions.

In differential geometry, the allowable morphisms are typically the *smooth maps*. A map $f : M^m \rightarrow N^n$ is *smooth* if for charts $U_i \subset M, U_j \subset N$, the composition $\phi_j \circ f \circ \phi_i^{-1}$ is a smooth map from the appropriate subset of \mathbb{R}^m to the appropriate subset of \mathbb{R}^n . That is, the co-ordinate maps have continuous partial derivatives of all orders. The *category of smooth manifolds*, denoted **DIFF** has as objects smooth manifolds, and as morphisms smooth maps. An invertible smooth map is called a *diffeomorphism*; the group of diffeomorphisms of a smooth manifold is typically a huge, unmanageable object, but certain features of it can be studied.

4.4. Topology. Basic notions of incidence or connectivity are part of our fundamental geometric intuition. Concepts such as “inside” and “outside”, or “bounded” and “unbounded” are topological. Topology can be thought of as the study of the *qualitative* properties of a space that are left unchanged under continuous deformations of the space; that is, deformations which may bend or stretch the space but do not cut or tear it. An allowable morphism between topological spaces is just a continuous map; invertible morphisms are called *homeomorphisms*.

4.4.1. Continuous maps. Topologists frequently discuss spaces far more abstract than manifolds. The concept of continuity in this general context relies on the definition of the following structure on a space.

Definition 4.16. A *topology* on a set X is a collection of subsets of X

$$\mathcal{U} \subset \{U \subset X\}$$

with the following properties:

- The empty set and X are both in \mathcal{U} .
- If U_1, \dots, U_n are a finite collection of elements in \mathcal{U} then $\cap_i U_i$ is in \mathcal{U} .
- If $\mathcal{V} \subset \mathcal{U}$ is an arbitrary collection of elements in \mathcal{U} , then $\cup_{V \in \mathcal{V}} V$ is in \mathcal{U} .

Thus, a topology is a system of subsets X which includes the empty set and X , and is closed under finite intersections and arbitrary unions. The sets in \mathcal{U} are called the *open* sets in X . Their complements are called the *closed* sets. From the definition, finite unions and arbitrary intersections of closed sets are closed.

Remark 4.17. It suffices, when defining a topology, to give a set of subsets of the space which are supposed to be open, and then let the open sets be the smallest collection of subsets, including the given sets, which satisfy the axioms of a topology. We call the given collection of sets a *basis* for the topology. Most of the spaces we will encounter have a *countable* basis.

Definition 4.18. Given a set $Y \subset X$, the *closure* of Y is the intersection of the closed subsets of X containing Y . Given a set Y , the *interior* of Y is the complement of the closure of the complement of Y . It is the union of all the open sets in X contained in Y .

Example 4.19. Let $Y \subset X$ be a subset. The *subspace* topology on Y is the topology whose open sets are the intersections $U \cap Y$ where U is open in X .

We will not discuss topological spaces in general in the sequel and stick only to some very concrete examples.

Let’s suppose we have two spaces X and Y where we understand what the open sets are. For instance, in \mathbb{R} the open sets are the unions of intervals of the kind (a, b) , where we don’t include the endpoints. In this context we can define abstractly what is meant by a continuous map.

Definition 4.20. A map from X to Y is *continuous* if the inverse of any open set is open.

Example 4.21. Let \sim be an equivalence relation on X , and let $\pi : X \rightarrow X/\sim$ be the quotient map to the space of equivalence classes. The *quotient* topology on X/\sim is the topology whose open sets are those $U \subset X/\sim$ such that $\pi^{-1}(U)$ is open in X . Thus, X/\sim has as many open sets as it is allowed subject to the condition that π is continuous.

Definition 4.22. A homeomorphism from X to Y is a continuous map which is invertible and has a continuous inverse.

The category of topological manifolds is denoted **TOP** and has as objects topological manifolds and as morphisms all continuous maps.

The group of homeomorphisms of a space are typically even larger and harder to understand than groups of diffeomorphisms. The advantage of working with arbitrary continuous maps is that many natural maps and constructions are on the face of them continuous rather than smooth, and one can accomplish more by allowing oneself greater flexibility.

We list some frequently encountered topological concepts:

Definition 4.23. A *neighborhood* of a point $p \in X$ is any open set $U \in X$ containing p .

Definition 4.24. A topological space is *Hausdorff* if for any two distinct points $p, q \in X$ there are neighborhoods of p and of q which are disjoint.

Definition 4.25. A manifold is defined much as a smooth manifold, with charts U_i and functions $\phi_i : U_i \rightarrow \mathbb{R}^n$ for some (typically fixed) n , but now we don't require that the transition functions $\phi_j \phi_i^{-1}$ be smooth, merely homeomorphisms. Formally, a manifold is a Hausdorff topological space with a countable basis, such that every point has a neighborhood homeomorphic to an open subset of \mathbb{R}^n for some (usually fixed) n .

Definition 4.26. A space X is *connected* if there are no proper nonempty subsets $U \subset X$ which are both *closed* and *open*. A space is *locally connected* if for any point p and any open set $O \subset X$ there is another $U \subset O$ such that U is connected, as a subspace of X .

Definition 4.27. A subset $X \subset Y$ (perhaps all of Y) is *compact* if it is closed, and for every collection \mathcal{U} of open sets in Y whose union contains X , there is a *finite subcollection* whose union contains X . A space X is *locally compact* if every $p \in X$ has a neighborhood whose closure is compact.

REFERENCES

- [1] R. Benedetti and C. Petronio, *Lectures on Hyperbolic Geometry*, Springer-Verlag Universitext (1992)
- [2] J. Conway, *The sensual quadratic form*, Math. Ass. Amer. Carus mathematical monographs **26**, (1997)
- [3] M. Gromov, *Metric structures for Riemannian and non-Riemannian spaces*, Birkhauser (1999)
- [4] R. Lickorish, *A finite set of generators for the homeotopy group of a 2-manifold*, Proc. Camb. Phil. Soc. **60** (1964), pp. 769–778
- [5] J. Montesinos, *Classical tessellations and three-manifolds*, Springer-Verlag (1987)
- [6] T. Rado, *Über den Begriff der Riemannsche Fläche*, Acta Univ. Szeged **2** (1924–26), pp. 101–121
- [7] J. Ratcliffe, *Foundations of hyperbolic manifolds*, Springer-Verlag GTM **149** (1994)
- [8] K. Schütte and B. L. van der Waerden, *Auf welcher Kugel haben 5, 6, 7, 8 oder 9 Punkte mit Mindestabstand Eins Platz*, Math. Ann. **123** (1951), pp. 96–124
- [9] J. Stallings, *Category language*, notes; available at <http://www.math.berkeley.edu/~stall>
- [10] W. Thurston, *Three-dimensional geometry and topology, vol. 1*, Princeton University Press, Princeton Math. Series 35 (1997)

DEPARTMENT OF MATHEMATICS, HARVARD, CAMBRIDGE, MA 02138
E-mail address: dannyca@math.harvard.edu