# ABSTRACT

DAVIDSON, RUTH ELLEN. Some Problems in Geometric Combinatorics and Mathematical Phylogenetics. (Under the direction of Patricia Hersh and Seth Sullivant.)

In this thesis, geometric objects essential to the underlying organization of problems in topological combinatorics and mathematical phylogenetics are studied. In particular, (1) lexicographic shellability characterizations of geometric and semimodular lattices are given, (2) results are given about a simplicial complex associated to a novel formulation of a hypergeometric identity using the Euler-Poincaré relation, (3) a complete combinatorial description is given of a family of polyhedral cones associated to the distance-based phylogenetic reconstruction algorithm UPGMA, and (4) polyhedral geometry is used to analyze the behavior of the distance-based phylogenetic methods UPGMA, least-squares phylogeny, and Neighbor-Joining near a polytomy.

Geometric lattices are characterized as those finite, atomic lattices such that every atom ordering induces a lexicographic shelling of the order complex given by an edge-labeling of the Hasse diagram known as a minimal labeling. Equivalently, geometric lattices are shown to be exactly those finite lattices such that every ordering on the join-irreducibles induces a lexicographic shelling. This new characterization fits into a similar paradigm as Peter McNamara's characterization of supersolvable lattices as those lattices admitting a different type of lexicographic shelling, namely one in which each maximal chain is labeled with a permutation of $\{1, \ldots, n\}$. A similar characterization of semimodular lattices is also given, and the relationship between two types of lexicographic shelling applied to finite graded atomic lattices is briefly explored.

A family of non-pure simplicial complexes $\{\Delta(n) : n \in \mathbb{Z}_{>0}\}$ is studied, of which the alternating sum of the face numbers is equal to one side of a known hypergeometric identity due to Alfred Dixon. The complex $\Delta(n)$ is shown to be shellable for all $n$. The Betti numbers of $\Delta(n)$ are calculated for $n \leq 4$, and the Euler-Poincaré relation is used to give a new proof of the identity for $n \leq 4$.

Distance-based phylogenetic algorithms attempt to solve the NP-hard least-squares phylogeny problem by mapping an arbitrary dissimilarity map on a set of taxa $X$ representing biological data to a tree metric on $X$. The input space of all dissimilarity maps on $X$ is identified as a Euclidean space which properly contains the set of all tree metric outputs on $X$ as a polyhedral fan. A distance-based phylogenetic method $f$ induces a partition of the input space into a family of regions

$$\{C(T) : f(x) \text{ is a tree metric realized by the combinatorial tree } T \text{ for all } x \in C(T)\}.$$

When the decision criteria for a method are linear, each region $C(T)$ is a polyhedral cone with an $\mathcal{H}$-representation given by the criteria. An instance of this type of family of polyhedral cones associated to the distance-based phylogenetic reconstruction algorithm UPGMA is studied. The set of extreme rays of a UPGMA cone $C(T)$ is shown to have a closed-form description in terms of the elements of the maximal chain in the lattice of set partitions $\Pi_n$ corresponding to a combinatorial tree $T$ on $n$ leaves. The spherical volumes of the UPGMA cones are computed for $n \leq 7$.

Phylogenetic inference on biological data routinely returns a tree with a *polytomy*, or an internal vertex with more than three neighbors, representing either a multi-way speciation event or a lack of sufficient data to resolve a binary phylogeny. Polyhedral geometry is used to compare the local nature of the subdivisions of the input space near a tree metric with a polytomy induced by the distance-based methods least-squares phylogeny, UPGMA, and Neighbor-Joining. The results of this investigation suggest that UPGMA and Neighbor-Joining poorly match least-squares phylogeny when the true tree has a polytomy with exactly four neighbors.

Some Problems in Geometric Combinatorics
and Mathematical Phylogenetics

by
Ruth Ellen Davidson

A dissertation submitted to the Graduate Faculty of
North Carolina State University
in partial fulfillment of the
requirements for the Degree of
Doctor of Philosophy

Mathematics

Raleigh, North Carolina

2014

APPROVED BY:

| | |
|---|---|
| Ezra Miller | Nathan Reading |

| | |
|---|---|
| Patricia Hersh | Seth Sullivant |
| Co-chair of Advisory Committee | Co-chair of Advisory Committee |

# DEDICATION

For Bernhard.

# BIOGRAPHY

Ruth Ellen Davidson was born in Duluth, Minnesota on October 30, 1976, to Donald William and Ruthanna Marie Davidson. She has an older sister, Charlotte Jeanne Karsh. She was a non-traditional student who did not discover she liked mathematics until she was 28 and took Calculus I as an elective at Seattle Central Community College. She is a composer, guitarist, cellist, vocalist, and bassist who has appeared on dozens of recordings in the genres of heavy metal, progressive rock, thrash, indie rock, and jazz. She is married to Bernhard John Kovacevich.

# ACKNOWLEDGEMENTS

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

# Chapter 1

# Introduction and Background Material

To paraphrase Phil Hanlon [44] *combinatorics* is the branch of mathematics that organizes discrete data in ways that make it manageable to analyze. We define *geometric combinatorics* as the sub-branch of combinatorics that either uses geometric objects to organize discrete data or studies geometric objects from a combinatorial point of view. In this thesis we study geometric objects essential to the underlying organization of problems in topological combinatorics and mathematical phylogenetics. In particular, two types of geometric objects that (1) play key roles in the results we obtain and (2) we study from a combinatorial viewpoint are simplicial complexes (Definition 1.1.1) and polyhedra (Definition 1.2.1). In a broader context, simplicial complexes are fundamental to the field of algebraic topology, and polyhedra are fundamental to Euclidean geometry.

Many key ideas in polyhedral geometry and simplicial topology begin with the work of Leonhard Euler (1707-1783). Euler's 1736 solution [33] to the famous Seven Bridges of Königsburg Problem is considered to be the first topological argument [20]. The problem was to find a path through the city of Königsburg using every one of its seven bridges; each of these seven bridges crossed the Pregel river exactly one time. The reason for interest in the problem was that it was impossible to show such a path did not exist by trial and error.

Euler is also responsible for the Euler polyhedron formula $v - e + f = 2$ for 3-dimensional polyhedra, where $v, e$, and $f$ denote the number of vertices, edges, and 2-dimensional faces, respectively. This formula has proven to be fundamental in algebraic topology as well as polyhedral geometry. The alternating sum of faces in arbitrary dimension is now known as the Euler characteristic and can be defined for either polyhedra or abstract simplicial complexes. For the specific applications in this thesis, we define the closely related *reduced* Euler characteristic for simplicial complexes in Definition 1.1.7.

So, the story of polyhedral geometry and simplicial topology begins in many ways with the work of Euler. The close relationship between the geometry and topology of polyhedra and simplicial complexes has persisted as the fields in which these objects are frequently used have developed. For example, Ludwig Schläfli (1814-1895) generalized the Euler polyhedron formula to convex polyhedra in arbitrary dimension [56] in 1852. In this particular monograph, Schläfli assumed that the boundary of a convex polytope was shellable [68]. Shellability, defined for simplicial complexes in Definition 1.3.2, is a concept that can be applied to both simplicial complexes and polyhedra, and gives similar information about them as geometric spaces. Shellable simplicial complexes are the main objects of interest in Chapters 2 and 3.

Another example of the historical relationship between polyhedral geometry and simplicial topology is found in the work of Henri Poincaré (1854-1912). Poincaré is considered to be largely responsible for combinatorial and algebraic topology developing into a mature field of study, and used polyhedra extensively in his work [20]. He obtained the Euler-Poincaré relation (Equation 1.1), which was published in 1899 [52]. This famous relation shows that the alternating sum of the Betti numbers (Definition 1.1.10) of a topological space is equal to the alternating sum of the number of cells forming the $i$-dimensional components of the space. Betti numbers can be viewed as numbers that count the $i$-dimensional holes of a space, so this relation combines two combinatorial perspectives on a geometric object. We exploit these perspectives in Chapter 3 to study a reformulation of a known identity due to Alfred Cardew Dixon (1865-1936) using a simplicial complex.

Once a mathematical object has been determined to be of interest, we may be able to better understand the object by studying a transformation of the object rather than directly studying the object itself. One well-known example of this strategy is the representation theory of groups, in which groups are studied via their images under homomorphisms into general linear groups of vector spaces. This philosophy is employed in topological combinatorics by associating a simplicial complex to a partially ordered set (Definition 1.1.2) called an order complex (Definition 1.1.3). The order complex encodes properties of the partially ordered set just as a group representation encodes information about a group. In particular, the types of shellings admitted by the order complex of a partially ordered set can reveal properties of the partially ordered sets. We examine this relationship in Chapter 2.

Another use of geometric combinatorics that is important in this thesis is the study of mathematical properties of the input and output spaces of an algorithm. Such mathematical properties may give insight about the performance of a method. For example, understanding the geometric structure or combinatorial properties of the input space of an algorithm can explain an observed bias of an algorithm on specific types of data sets. These insights can be exploited to create better algorithms and improve existing algorithms. When an algorithm takes inputs in a Euclidean space and makes decisions according to linear selection criteria, we can study

the decomposition of the input space indexed by the possible outputs of the algorithm using polyhedral geometry.

There are many tools available to aid in the study of such a decomposition found in the field of linear optimization. Linear optimization is the extensively developed area of mathematics concerned with minimizing or maximizing a function subject to a set of constraints, where the function and constraints are determined by linear equations. This type of optimization was pioneered in 1947 when George Dantzig designed a method for solving linear formulations of United States Air Force planning problems. The method he designed is now known as the simplex method and was soon determined to be applicable to a wide range of problems [22]. The fundamental objects organizing all linear optimization problems are polyhedra, which have been studied for centuries. However, the emergence of the field of linear optimization was facilitated by the development of high-performance computing due to the high dimension of the polyhedra often associated with linear optimization problems of practical interest.

The algorithms we will study using this approach, that is, the approach of using polyhedral geometry to study input and output spaces, are phylogenetic reconstruction algorithms. A *phylogeny* is a mathematical model of the common evolutionary history of a group of genes or a group of species, and is commonly represented using a phylogenetic tree (Definition 1.4.1). Sometimes the genes or species are referred to as taxonomical units or taxa. The phylogenetic reconstruction methods we will study use data collected under the assumption of models of DNA, RNA, or protein sequence evolution on the taxa. We assume there is a *phylogenetic signal* of the true history that the method is detecting in the data.

Distance-based phylogenetic reconstruction methods, introduced in Section 1.4.1, take a point in Euclidean space as an input and return a phylogenetic tree with edge weights as an output. These methods include polynomial-time algorithms that are relevant in the age of huge data sets. Two fundamental examples of these algorithms include the Neighbor-Joining (NJ) algorithm of Naruya Saitou and Masatoshi Nei (Algorithm 5.4.6) and the Unweighted Pair Group Method with Arithmetic Mean (UPGMA) of Robert Sokal and Charles Michener (Algorithm 4.2.1). Both NJ and UPGMA use linear selection criteria and therefore define families of polyhedra that are indexed by the combinatorial type of tree returned by the algorithm.

It is the within the purview of biologists to develop heuristics for solving computationally difficult problems, and it is within the purview of mathematicians to analyze how well these heuristics perform by posing and studying mathematically tractable questions associated to the heuristics. NJ and UPGMA are considered by biologists to have certain drawbacks in terms of accuracy. For example, they do not always satisfy performance criteria such as the ability to produce the correct tree on a fixed data set, where the correct tree is known because the data was produced by simulations or curated from a set of taxa with a known evolutionary history. Also, these algorithms may exhibit bias that contradicts the expected values of model

distributions on phylogenetic trees. Using the mathematical properties of the input and output spaces associated to an algorithm we may confirm or further explain biases and pathologies already known to biologists by other means.

There exist phylogenetic reconstruction methods that are considered superior by biologists due to the use of more accurate models of evolution or superior performance on real data sets, such as maximum likelihood estimation. However, the input and output spaces of methods that are not distance-based are often difficult to study mathematically. There are also distance-based methods that include corrections and improvements to the approach used in NJ that include the Weighted Neighbor-Joining method "Weighbor" of William Bruno, Nicholas Socci, and Aaron Halpern [21] and the BIONJ algorithm of Olivier Gascuel [38]. But both Weighbor and BIONJ induce partitions of the input space that are not polyhedral, and so are not amenable to the powerful toolkits at hand from linear optimization.

So in some ways, NJ and UPGMA should be viewed as fundamental examples in a large class of algorithms. By studying the input and output spaces for NJ and UPGMA using polyhedral geometry, we hope to gain insight into the issues observed by biologists. Studying these algorithms may suggest what can be learned about other methods. Furthermore, these algorithms give rise to families of polyhedra that are inherently interesting from the perspective of geometric combinatorics. In Chapter 4, we explore the combinatorial properties of the family of polyhedra associated to UPGMA. The organizational tool we use to explain these combinatorial properties is a partially ordered set (Definition 1.1.2) known as the lattice of set partitions $\Pi_n$ (Example 1.1.4). The partially ordered set $\Pi_n$ can be studied using the concept of shellability of an order complex (Definition 1.1.3). This concept is fundamental to the results in Chapter 2, and provides another link between simplicial complexes and polyhedral geometry.

The NJ and UPGMA algorithms can be viewed as polynomial-time heuristics for the NP-hard least-squares phylogeny (LSP) problem (Problem 5.1.1). This provides a concrete starting point for the use of our geometric perspective in Chapter 5, in which we examine the decomposition of the input space in certain regions corresponding to data that may not be sufficient to resolve a complete evolutionary history. These regions correspond to phylogenies that model non-binary speciation events, or trees with an unresolved branching structure. This type of tree contains a high-degree internal vertex called a polytomy (Definition 1.4.3). Trees containing polytomies are routinely returned by many types of phylogenetic reconstruction methods, not just the methods addressed in this thesis.

Chapter 1 is devoted to developing background and terminology necessary to present the topics in later chapters. Section 1.1 explains the connection between simplicial complexes and partially ordered sets that is fundamental to Chapter 2, as well as key concepts from simplicial topology fundamental to Chapter 3. Section 1.2 explains some basic concepts in polyhedral geometry that are used in Chapters 3 and 4. Section 1.3 defines shellability and shellable simplicial

complexes, and briefly explores the notions of lexicographic shellability used in Chapter 2 and general non-pure shellability used in Chapter 3. Section 1.4 lays out the fundamental problem of phylogenetic inference, explains the notion of distance-based phylogenetic inference that is examined in detail in Chapters 4 and 5, and sets up a geometric perspective on distance-based phylogenetic inference using the ideas from polyhedral geometry outlined in Section 1.2.

The results in Chapter 2 about lexicographic shellability characterizations of geometric and semimodular lattices are joint work with Patricia Hersh and appear in the paper [23]. The results in Chapters 4 and 5 are joint work with Seth Sullivant and appear in the papers [24] and [26], respectively.

## 1.1   Simplicial Complexes and Partially Ordered Sets

The first object from geometric combinatorics that we will introduce is the simplicial complex. Simplicial complexes can be thought of in many ways, including topological spaces, geometric objects, or set systems for organizing information. We will primarily be concerned with their topological properties and their relationship to the theory of partially ordered sets.

**Definition 1.1.1.** An *(abstract) simplicial complex* on a vertex set $V$ is a collection $\Delta$ of subsets of $V$ satisfying

1. if $v \in V$ then $\{v\} \in \Delta$, and

2. if $F \in \Delta$ and $G \subseteq F$, then $G \in \Delta$.

The subsets of $V$ comprising $\Delta$ are called *faces* or *simplices*. The dimension $\dim F$ of a face $F$ is $|F| - 1$, and $\dim \Delta$ is simply $\max\{\dim F : F \in \Delta\}$. A face $F$ is a *facet* if $F$ is not properly contained in any other face of $\Delta$. We say $\Delta$ is *pure* if all the facets of $\Delta$ have the same dimension. We write $\overline{F}$ to denote the sub-complex of $\Delta$ generated by $F$, or in other words $\overline{F} = \{G \in \Delta : G \subseteq F\}$.

Figure 1.1 shows two simplicial complexes on the vertex set $V = \{1, 2, 3, 4, 5\}$. In Figure 1.1-(a) we have a pure simplicial complex

$$\Delta_a = \{\emptyset, \{1\}, \{2\}, \{3\}, \{4\}, \{5\}, \{1,2\}, \{1,4\}, \{2,3\}, \{2,4\}, \{2,5\}, \{3,5\}, \{1,2,4\}, \{2,3,5\}\}$$

and in Figure 1.1-(b) we have a non-pure simplicial complex

$$\Delta_b = \{\emptyset, \{1\}, \{2\}, \{3\}, \{4\}, \{5\}, \{1,2\}, \{1,4\}, \{2,3\}, \{2,4\}, \{2,5\}, \{3,5\}, \{1,2,4\}\}.$$

Both $\Delta_a$ and $\Delta_b$ contain the face $F = \{1, 2, 4\}$, and so

$$\overline{F} = \{\emptyset, \{1\}, \{2\}, \{4\}, \{1, 2\}, \{1, 4\}, \{2, 4\}, \{1, 2, 4\}\}.$$

is a sub-complex of both $\Delta_a$ and $\Delta_b$.



Figure 1.1: Two simplicial complexes such that $V = \{1, 2, 3, 4, 5\}$.

Partially ordered sets, commonly referred to as *posets*, play key roles throughout combinatorics via topics such as Möbius inversion, hyperplane arrangements, and applications to group theory. Posets are useful in computer science and computational biology, and are applied in various other branches of mathematics. We direct the interested reader to Chapter 3 of Richard Stanley's *Enumerate Combinatorics Volume I* [66] for a thorough introduction to the theory of posets.

**Definition 1.1.2.** A *partially ordered set* or *poset* is a set $P$ with a binary order relation satisfying three axioms:

1. For all $p \in P$, $p \le p$. (reflexivity)

2. If $p, q \in P$ satisfy both $p \le q$ and $q \le p$, then $p = q$. (antisymmetry)

3. If $p, q, r \in P$ satisfy both $p \le q$ and $q \le r$, then $p \le r$. (transitivity)

Note: throughout this thesis all posets will be taken to be finite.

Two elements $p$ and $q$ of a poset $P$ are *comparable* if either $p \le q$ or $q \le p$ holds; else we say $p$ and $q$ are *incomparable*. The poset $P$ *has a* $\hat{0}$ if there exists an element $\hat{0}$ such that $\hat{0} \le p$ for all $p \in P$, and $P$ *has a* $\hat{1}$ if there exists an element $\hat{1}$ such that $\hat{1} \ge p$ for all $p \in P$. A *chain* $C$ is a poset that is *totally ordered* or admits a *linear order*, i.e. every pair $\{p, q\} \subset C$ is comparable. Examples of totally ordered sets include the real numbers $\mathbb{R}$ and the set of the first $n$ positive integers $\{1, 2, \ldots, n\}$, denoted $[n]$.

A subset $C$ of a poset $P$ is a *chain of $P$* if $C$ is a chain when regarded as a subposet of $P$. If there does not exist a larger chain of $P$ containing $C$, we say that $C$ is *maximal*. If there does not exist $p, q \in C$ and $r \in P \setminus C$ satisfying (i) $p < r < q$ and (ii) $C \cup \{r\}$ is a chain of $P$, then we say $C$ is *unrefinable* or *saturated*. If $p < q$ and there does not exist $r \in P$ such that $p < r < q$, we say that $q$ *covers* $p$ and write $p \lessdot q$.

To visualize and study a poset $P$ we associate the useful *Hasse diagram* of $P$ which is simply a graph with vertices given by the elements of $P$ and edge set $\mathcal{E}(P)$ given by the covering relations of $P$. We follow the convention that when $p \lessdot q$, $q$ is drawn above $p$. The Hasse diagram of a poset is shown in Figure 1.2-(a).



Figure 1.2: A Hasse diagram and an order complex of a poset.

7

### 1.1.1 Order Complexes of Partially Ordered Sets

Next we introduce a family of simplicial complexes associated to posets known as order complexes. Order complexes of posets are useful because we can obtain valuable information about the poset by studying the associated order complex. Such information may include determination of membership in an important class of posets, as we will see in Chapter 2.

**Definition 1.1.3.** Let $P$ be a poset. The *order complex* $\Delta(P)$ of $P$ is the simplicial complex with vertex set given by the elements of $P$ and faces given by the chains of $P$.

Figure 1.2-(b) shows the order complex of the poset whose Hasse diagram is shown in Figure 1.2-(a). Note that the maximal chains with 3 elements correspond to 2-simplices (also known as triangles) in the order complex. The non-maximal chain $\{\hat{0}, \hat{1}\}$ corresponds to the 1-simplex (also known as an edge) in the order complex that is the intersection of the 3 maximal 2-simplices in the order complex.

Another partially ordered set that appears many times in this thesis is the lattice of set partitions, which we introduce in the next example. Lattices are a special type of poset; the precise definition of a lattice is given in Definition 2.2.1.

**Example 1.1.4.** Let $\Pi_n$ consist of all set partitions of a set with $n$ elements. Without loss of generality we can identify this underlying set as $[n]$. We consider the set $[n]$ and the blocks in a set partition of $[n]$ as consisting of unordered elements: $\lambda_1 | \ldots | \lambda_k$ denotes a set partition with $k$ blocks where $\lambda_i \subset [n]$ and $\lambda_i \cap \lambda_j = \emptyset$ for $i \neq j$, $i, j \in [k]$. When the context is clear, we will use (for example) $12|345$ as shorthand for the set partition $\{\{1, 2\}, \{3, 4, 5\}\}$ of $[5]$.

Set partitions in $\Pi_n$ are ordered by *refinement*, which means that $\lambda_1 | \ldots | \lambda_k \leq \mu_1 | \ldots | \mu_\ell$ if and only if for each $i \in [k]$ there exists a $j \in [\ell]$ satisfying $\lambda_i \subseteq \mu_j$.

The Hasse diagram of $\Pi_3$ is shown in Figure 1.3. Two posets $P$ and $Q$ are *isomorphic* if there exists an order-preserving bijection $\phi : P \to Q$ whose inverse is order-preserving, that is $s \leq t$ in $P$ if and only if $\phi(s) \leq \phi(t)$ in $Q$. Note that $\Pi_3$ is isomorphic to the poset shown in Figure 1.2-(a).

### 1.1.2 Some Topological and Combinatorial Invariants of Simplicial Complexes

One way to understand the structure of a simplicial complex is by examining its topological and combinatorial qualities. Any geometric realization of an abstract simplicial complex is a topological space, and we can understand this space combinatorially. One useful combinatorial invariant simply counts the faces of each dimension of a finite simplicial complex $\Delta$:

Figure 1.3: $\Pi_3$, the lattice of set partitions of the set $\{1, 2, 3\}$.

**Definition 1.1.5.** The *f-vector* $f_\Delta = (f_0, f_1, \ldots, f_d)$ is the integer vector with entries $f_i$ counting the number of faces of dimension $i$. The maximal entry $f_d$ counts the number of facets of $\Delta$, and $\dim \Delta = d$. If we consider the empty set to be a face of a simplicial complex $\Delta$, we say $\emptyset$ is a face with dimension equal to $-1$, and $f_\Delta = (f_{-1}, f_0, f_1, \ldots, f_d)$, where $f_{-1} = 1$.

**Example 1.1.6.** When we defined $\Delta_a$ and $\Delta_b$, shown in Figure 1.1-(a) and Figure 1.1-(b), we specified that $\emptyset$ was a face of both $\Delta_a$ and $\Delta_b$. So $f_{\Delta_a} = (1, 5, 6, 2)$ and $f_{\Delta_b} = (1, 5, 6, 1)$. If $\overline{\Delta_a}$ is the simplicial complex where we *do not* consider $\emptyset$ to be a face, then $f_{\overline{\Delta_a}} = (5, 6, 2)$.

Another combinatorial invariant arises as the alternating sum of the entries in the $f$-vector $f_\Delta$:

**Definition 1.1.7.** The *reduced Euler characteristic* of the simplicial complex $\Delta$ is the alternating sum

$$\widetilde{\chi}(\Delta) = \sum_{i=-1}^{d} (-1)^i f_i$$

where $f_\Delta = (f_{-1}, f_0, \ldots, f_d)$.

Note that the alternating sum in Definition 1.1.7 of the $f$-vector where $\emptyset$ is *not* included as a face is simply called the *Euler characteristic*. So, the modifier *reduced* in this context specifically indicates the inclusion of $\emptyset$ as a face.

Topological qualities of a simplicial complex $\Delta$ can determine the reduced Euler characteristic of $\Delta$, as the next example demonstrates. Recall that a *homotopy* between two functions $f$ and $g$ from a space $X$ to a space $Y$ is a continuous map $G$ from $X \times [0,1] \to Y$ such that $G(x,0) = f(x)$ and $G(x,1) = g(x)$.

**Definition 1.1.8.** Two topological spaces $X$ and $Y$ are *homotopy equivalent* if there exist continuous maps $j : X \to Y$ and $k : Y \to X$ such that $j \circ k$ is homotopic to the identity map on $Y$ and $k \circ j$ is homotopic to the identity map on $X$.

A useful way to think of homotopy equivalence of two spaces is that one can be continuously deformed into the other. So the capital letters "A" and "O" are homotopy equivalent.

**Example 1.1.9.** Note that if $\Delta = \{\emptyset, \{1\}\}$, then $\widetilde{\chi}(\Delta) = (-1)^{-1} \times (1) + (-1)^0 \times (1) = 0$. Also, consider again the complex $\Delta_a$ from Figure 1.1-(a). We have

$$\widetilde{\chi}(\Delta_a) = (-1)^{-1} \times (1) + (-1)^0 \times (5) + (-1)^1 \times (6) + (-1)^2 \times (2) = 0$$

The complexes in Example 1.1.9 are instances of the fact that the reduced Euler characteristic of any simplicial complex that can be continually deformed to a point is 0. The complexes $\Delta$ and $\Delta_a$ are equivalent as topological spaces under homotopy equivalence. However, it is not the case that every simplicial complex with a reduced Euler characteristic of 0 is homotopy equivalent to a point, as we will see in Chapter 3.

A very useful set of topological invariants of a simplicial complex is the set of Betti numbers. To define Betti numbers we must first understand the notion of the (simplicial) homology groups of a simplicial complex. For an introduction to simplicial homology, see for example Section 2.1 of *Algebraic Topology* by Allen Hatcher [42]. For a simplicial complex $\Delta$, let $\Delta_k$ denote the set of all $k$-dimensional simplices in $\Delta$, i.e. the set of all simplices in $\Delta$ with $k+1$ vertices.

A *simplicial k-chain* is a formal sum of $k$-simplices $\sum_{i=1}^{j} c_i \sigma_i$ where $\sigma_i \in \Delta_k$ and $c_i \in \mathbb{Z}$. Note that we will use the term "chain" in this context to mean something completely different than its usage in conjunction with posets as defined earlier. Let $C_k$ denote the free abelian group with the basis given by the elements of $\Delta_k$. The group $C_k$ is often called a *chain group*. Let $\sigma = \{v_1, \ldots, v_{k+1}\} \in \Delta_k$. The *kth boundary map* $\partial_k : C_k \to C_{k-1}$ between chain groups is the function defined by

$$\partial_k(\sigma) = \sum_{m=1}^{k+1} (-1)^m \{v_1, \ldots, \widehat{v_m}, \ldots, v_{k+1}\}$$

where $\{v_1, \ldots, \widehat{v_m}, \ldots, v_{k+1}\}$ is the $(k-1)$-simplex obtained by omitting the vertex $v_m$. The elements of the subgroup $\ker \partial_k$ of $C_k$ are called *cycles* and the elements of the subgroup $\operatorname{im} \partial_{k+1}$

of $C_k$ are called *boundaries*. It is simple to verify that $\operatorname{im}\partial_{k+1} \subset \ker\partial_k$, so that the quotient group $H_k = \ker\partial_k / \operatorname{im}\partial_{k+1}$ is defined. We call the group $H_k$ the *kth homology group* of $\Delta$, and also write $H_k(\Delta)$ when the specific simplicial complex under discussion must be made clear.

**Definition 1.1.10.** The number $\beta_k(\Delta) = \operatorname{rank}(H_k(\Delta))$ is the *kth Betti number* of $\Delta$.

The number $\beta_k(\Delta)$ can often be regarded as the number of $k$-dimensional holes that $\Delta$ has as a topological space. When it is clear from context which simplicial complex we are discussing, we simply write $\beta_k$. We have $\beta_{-1} = 1$, for the single generator $\emptyset$ of the chain group $C_{-1}$, as all chain groups $C_k$ for $k < -1$ are zero. The number $\beta_0$ counts the number of connected components of $\Delta$, and the number $\beta_1$ counts the number of 1-dimensional holes.

If $d$ is the maximum dimension of a face of a simplicial complex $\Delta$, the face numbers and Betti numbers of $\Delta$ are related by the Euler-Poincaré relation, which is attributed [9] to Henri Poincaré:

$$\sum_{i=-1}^{d} (-1)^i f_i = \sum_{i=-1}^{d} (-1)^i \beta_i. \tag{1.1}$$

**Example 1.1.11.** We check Equation 1.1 for $\Delta_a$ and $\Delta_b$ from Figure 1.1. We computed $\widetilde{\chi}(\Delta_a) = 0$ in Example 1.1.9. It is clear that $\Delta_a$ has one connected component, so $\beta_0(\Delta_a) = 1$, and no higher-dimensional holes, so $\beta_k(\Delta_a) = 0$ for $k > 0$. Thus,

$$(-1)^{-1} \times \beta_{-1}(\Delta_a) + (-1)^0 \times \beta_0(\Delta_a) = -1 + 1 = 0.$$

So Equation 1.1 holds for $\Delta_a$.

We have

$$\widetilde{\chi}(\Delta_b) = (-1)^{-1} \times (1) + (-1)^0 \times (5) + (-1)^1 \times (6) + (-1)^2 \times (1) = -1.$$

To compute the Betti numbers of $\Delta_b$, we note that $\Delta_b$ has one connected component, so that $\beta_0(\Delta_b) = 1$, but also one 1-dimensional hole made by the triangle with edges $\{2,3\}, \{2,5\}$ and $\{3,5\}$. So, $\beta_1(\Delta_b) = 1$. Therefore $\sum_{i=-1}^{2}(-1)^i\beta_i(\Delta_b) = -1$, and Equation 1.1 also holds for $\Delta_b$.

## 1.2 Polyhedral Geometry

Polyhedra are fundamental objects of study in Euclidean geometry. In particular, they are used as tools in linear optimization. See Günter Ziegler's *Lectures on Polytopes* [72] for a comprehensive theoretical introduction to polyhedral geometry. For an introduction to polyhedral geometry from the perspective of linear optimization, we recommend either *Introduction to*

*Linear Optimization* by Dimitris Bertsimas and John Tsitsiklis [7] or Vašek Chvátal's *Linear Programming* [22].

We now develop some terminology essential to polyhedral geometry in a Euclidean space $\mathbb{R}^d$. Given a subset $X \subseteq \mathbb{R}^d$, the *convex hull* of $X$, denoted conv$(X)$, is the set of all convex combinations of points in $X$, or in other words:

$$\text{conv}(X) = \left\{ \lambda_1 \mathbf{x}_1 + \cdots + \lambda_k \mathbf{x}_k : \{\mathbf{x}_1, \ldots, \mathbf{x}_k\} \subseteq X, \ \lambda_i \geq 0, \ \sum_{i=1}^{k} \lambda_i = 1 \right\}.$$

The *conical hull* of $X$ is the set of all nonnegative combinations of points in $X$:

$$\text{cone}(X) = \{ \lambda_1 \mathbf{x}_1 + \cdots + \lambda_k \mathbf{x}_k : \{\mathbf{x}_1, \ldots, \mathbf{x}_k\} \subseteq X, \ \lambda_i \geq 0 \}.$$

Let $\mathbf{c} \in \mathbb{R}^d$ and $z \in \mathbb{R}$. The (possibly empty) *half-space* in $\mathbb{R}^d$ defined by $\mathbf{c}$ and $z$ is the set of all $\mathbf{x} \in \mathbb{R}^d$ such that $\mathbf{c} \cdot \mathbf{x} \leq z$.

**Definition 1.2.1.** A $\mathcal{H}$-polyhedron $P$ is the intersection of finitely many half-spaces in $\mathbb{R}^d$.

If $P$ is an $\mathcal{H}$-polyhedron we may write $P = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{z}\}$, and say $A\mathbf{x} \leq \mathbf{z}$ is an *$\mathcal{H}$-representation of $P$*, where $A \in \mathbb{R}^{m \times d}$ and $\mathbf{z} \in \mathbb{R}^m$. A linear inequality $\mathbf{c}\mathbf{x} \leq c_0$ is *valid* for $P$ if it is satisfied for all points $\mathbf{x} \in P$. A *face* of $P$ is any set of the form

$$F = P \cap \{\mathbf{x} \in \mathbb{R}^d : \mathbf{c}\mathbf{x} = c_0\}$$

where $\mathbf{c}\mathbf{x} \leq c_0$ is a valid inequality for $P$. The *dimension* of a face is the dimension of its affine hull.

**Definition 1.2.2.** A *$\mathcal{V}$-polyhedron* is a subset of $\mathbb{R}^d$ that can be written as conv$(X)$ + cone$(Y)$ for finite subsets $X$ and $Y$ of $\mathbb{R}^d$.

In Definition 1.2.2, the symbol $+$ denotes the *Minkowski sum*: if $X, Y \subseteq \mathbb{R}^d$ then $X + Y = \{\mathbf{x} + \mathbf{y} : \mathbf{x} \in X, \mathbf{y} \in Y\}$. If $P = \text{conv}(X) + \text{cone}(Y)$ then conv$(X)$ + cone$(Y)$ is a *$\mathcal{V}$-representation* of $P$.

Theorem 1.2.3, due to Gyula Farkas (1847-1930) [35], is sometimes called "Farkas' Theorem" in the literature, and is not to be confused with the also well-known "Farkas' Lemma," to which Section 1.4 of Ziegler's book [72] is devoted. While we will not have need of Farkas' famous lemma in this thesis, we will make use of Theorem 1.2.3.

**Theorem 1.2.3.** *A subset $P \subseteq \mathbb{R}^d$ is an $\mathcal{H}$-polyhedron if and only if $P$ is a $\mathcal{V}$-polyhedron.*

Note that Theorem 1.2.3 merely asserts that if there *exists* an $\mathcal{H}$-representation of a polyhedron $P$ then $P$ also has a $\mathcal{V}$-representation, and vice versa. It is still a general problem of

Figure 1.4: A polyhedron in $\mathbb{R}^2$.

interest in polyhedral geometry to compute the other representation when one is known. This problem is often computationally intractable for "large" (in other words, useful) dimension $d$.

**Example 1.2.4.** Figure 1.4 shows a polyhedron $P \subseteq \mathbb{R}^2$. The dashed lines are valid inequalities for $P$ that give an $\mathcal{H}$-representation for $P$. The arrows show which side of the inequalities we identify as the half-space defined by each inequality. We fill in the dashed lines where they define a 1-dimensional face of $P$. The shaded grey area that is the intersection of the three half-spaces is $P$. The points $A, B$, and $C$ are the 0-dimensional faces of $P$, also known as vertices. The $\mathcal{V}$-representation of $P$ is $P = \mathrm{conv}\{A, B, C\}$. Note that $P$ is also a geometric realization of the abstract 2-simplex with vertex set $\{A, B, C\}$.

Some simpler versions of Theorem 1.2.3 will be of use to us in this thesis. In particular, we will have need of convenient descriptions of polyhedral cones and polytopes. We say a set $C \subseteq \mathbb{R}^d$ is a *polyhedral cone* if $C = \mathrm{cone}(Y)$ for a finite set of vectors $Y \subset \mathbb{R}^d$. Let $\mathbf{0}$ denote the zero vector in $\mathbb{R}^d$.

**Theorem 1.2.5.** *A set $C \subseteq \mathbb{R}^d$ is a polyhedral cone $C = \mathrm{cone}(Y)$ for some finite subset $Y \subseteq \mathbb{R}^d$ if and only if $C = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{0}\}$ for some $A \in \mathbb{R}^{m \times d}$.*

Hereafter we take *cone* to mean polyhedral cone.

13

Figure 1.5: A 2-dimensional polyhedral fan in $\mathbb{R}^3_{\geq 0}$.

**Definition 1.2.6.** An *extreme ray* **r** of a cone $C$ is a 1-dimensional face of $C$.

Every polyhedral cone can be written as the cone of its extreme rays, so in practice when we write $C = \text{cone}(Y)$ we may assume, if it is clear from the context, that $Y$ is precisely the set of vectors that generate extreme rays of $C$. We commonly refer to a point vector $\mathbf{v} \in \mathbb{R}^d$ that generates an extreme ray **r** interchangeably with the extreme ray **r** itself.

**Definition 1.2.7.** For any convex subset $P \subseteq \mathbb{R}^d$, the *lineality space* of $P$ is defined as

$$\text{lineal}(P) = \{\mathbf{y} \in \mathbb{R}^d \; : \; \mathbf{x} + t\mathbf{y} \in P \;\; \text{for} \;\; \text{all} \; \mathbf{x} \in P, t \in \mathbb{R}\}$$

When a polyhedron $P$ satisfies $\text{lineal}(P) = \{\mathbf{0}\}$, then we say $P$ is *pointed*. In Chapters 4 and 5 we will use the idea of a polyhedral fan:

**Definition 1.2.8.** A *fan* is a family $\mathcal{F}$ of cones in $\mathbb{R}^d$ such that:

1. if $P \in \mathcal{F}$ then every nonempty face of $P$ is in $\mathcal{F}$, and

2. if $P_1, P_2 \in \mathcal{F}$ then $P_1 \cap P_2$ is a face of both $P_1$ and $P_2$.

Informally, a polyhedral fan is a family of cones that is easy to work with because cones in the family fit together nicely. The *dimension* of a fan $\mathcal{F}$ is the largest dimension of a cone in $\mathcal{F}$. Figure 1.5 shows a 2-dimensional fan in the *positive octant*, also known as $\mathbb{R}^3_{\geq 0}$, of $\mathbb{R}^3$. The labeled points in the picture generate the extreme rays of the fan. Each 2-dimensional cone is the cone of two extreme rays.

The next version of Farkas' Theorem concerns polytopes. A polyhedron is *bounded* if it contains no ray; a *polytope* is a bounded polyhedron. A $\mathcal{V}$-polytope is the convex hull of a finite

point set in $\mathbb{R}^d$. An $\mathcal{H}$-polytope is a bounded $\mathcal{H}$-polyhedron. The polyhedron in Figure 1.4 is a polytope.

**Theorem 1.2.9.** *A subset $P \subseteq \mathbb{R}^d$ is a $\mathcal{V}$-polytope if and only if it is an $\mathcal{H}$-polytope.*

The concept of a polytopal complex is similar to the concept of a polyhedral fan. We introduce it now to motivate the idea of shellability that is introduced in Section 1.3.

**Definition 1.2.10.** A *polytopal complex* is a finite, nonempty collection $\mathcal{C}$ of polytopes, (called the *faces* of $\mathcal{C}$) in $\mathbb{R}^d$ that contains all the faces of its polytopes, and such that the intersection of two polytopes $P_1, P_2 \in \mathcal{C}$ is a face of both $P_1$ and $P_2$. The *dimension* of a polytopal complex is the dimension of the largest polytope in $\mathcal{C}$, and $\mathcal{C}$ is *pure* if all the inclusion-maximal faces, or *facets* of $\mathcal{C}$, have the same dimension.

Note that, for example, the meanings of the terms *facet* and *dimension* for a polytopal complex are analogous to the meanings for these terms in regard to simplicial complexes (Definition 1.1.1).

**Example 1.2.11.** Given a polytope $P$, the *boundary complex* of $P$, denoted $\mathcal{C}(\partial P)$, is the polytopal complex formed by taking the collection of all proper faces of $P$. Then the facets of $\mathcal{C}(\partial P)$ are the facets of $P$, and $\mathcal{C}(\partial P)$ is a pure (dim $P$-1)-dimensional polytopal complex.

## 1.3   Shellability

The notion of shellability originated in polyhedral theory via the study of boundary complexes of convex polytopes, which are introduced in Example 1.2.11. As Michelle Wachs discusses in Lecture 3 of [68], the geometer Ludwig Schläfli (1814-1895) assumed the following theorem without proof in his 1852 manuscript *Theorie der vielfachen Kontinuität* [56].

**Theorem 1.3.1.** *The boundary complex of a convex polytope is shellable.*

Shläfli assumed the content of Theorem 1.3.1 in the course of computing the Euler characteristic of a convex polytope. Heinz Bruggesser and Peter Mani proved Theorem 1.3.1 in 1970 [17], and Peter McMullen used Theorem 1.3.1 in his proof of the upper bound conjecture for simplicial polytopes [48]. Since simplicial complexes arise outside of polyhedral geometry, shellability has many applications outside of polyhedral geometry as well. See Definition 8.1 in [72] for the definition of shellability for polytopal complexes. We will only define shellability for simplicial complexes.

**Definition 1.3.2.** A simplicial complex $\Delta$ is *shellable* if its facets can be arranged in a linear order $F_1, F_2, \ldots, F_t$ so that the subcomplex $\left( \bigcup_{i=1}^{k-1} \overline{F_i} \right) \cap \overline{F_k}$ is pure and (dim $F_k$−1)- dimensional for $k = 2, \ldots, t$. Such an ordering is called a *shelling*.

The *wedge* of two disjoint topological spaces $X$ and $Y$ is the quotient space obtained by identifying two points $x_0 \in X$ and $y_0 \in Y$ as equivalent. The next theorem, due to Anders Björner and Michelle Wachs, appears for general, i.e. not necessarily pure, simplicial complexes, in [15], and explains why shellable simplicial complexes have attractive topological properties.

**Theorem 1.3.3.** *A shellable simplicial complex has the homotopy type of a wedge of spheres in varying dimensions. For each dimension $r$, the number of $r$-spheres is the number of $r$-facets whose entire boundary is contained in the union of earlier facets in the shelling order.*

**Example 1.3.4.** The simplicial complex $\Delta_a$ in Figure 1.1-(a) is not shellable; there are two facets, $\{1, 2, 4\}$ and $\{2, 3, 5\}$, and no matter which facet we choose to come first in a proposed shelling order, $\overline{F}_1 \cap \overline{F}_2 = \{2\}$ which is a 0- dimensional sub-complex, whereas $(\dim F_2 - 1) = (\dim F_1 - 1) = 1$. However, the non-pure simplicial complex $\Delta_b$ in Figure 1.1-(b) is shellable. We can take $F_1 = \{1, 2, 4\}$, $F_2 = \{2, 3\}, F_3 = \{2, 5\}$, and $F_4 = \{3, 5\}$. It is easy to check that $\left( \bigcup_{i=1}^{k-1} \overline{F}_i \right) \cap \overline{F}_k$ is pure and $(\dim F_k - 1)$- dimensional for $k = 2, 3$, and 4 in this case.

### 1.3.1 Lexicographic Shellability

Let $P$ and $\Lambda$ be posets. Given a map $\lambda$ from the edges of the Hasse diagram of $P$ to $\Lambda$ $\lambda : \mathcal{E}(P) \to \Lambda$ we can associate a label sequence $\lambda(x_1, x_2), \lambda(x_2, x_3), \ldots, \lambda(x_{n-1}, x_n)$ to each saturated chain $C = x_1 \lessdot x_2 \lessdot \cdots \lessdot x_n$ in $P$. Different types of lexicographic shellabillity are defined by conditions on these label sequences arising from various edge-labelings. For all discussions in this thesis, it is sufficient to take $\Lambda$ to be the set $\mathbb{Z}$ of integers with the usual total ordering. We say that $C$ is *increasing* if the label sequence of $C$ is strictly increasing. The first type of lexicographic shellability that we examine is EL-shellability.

**Definition 1.3.5.** A map from the edges of the Hasse diagram of a poset $P$ to a poset $\Lambda$ $\lambda : \mathcal{E}(P) \to \Lambda$ is an *EL-labeling* for $P$ if

1. for every interval $[x, y]$ of $P$, there is a unique *rising chain* $C := x = x_1 \lessdot x_2 \lessdot \cdots \lessdot x_j = y$ where $\lambda(x, x_2) \leq \lambda(x_2, x_3) \leq \cdots \leq \lambda(x_{j-1}, y)$, and

2. the label sequence of $C$ is lexicographically smaller than the label sequence of every other saturated chain in the interval $[x, y]$.

Figure 1.6-(a) shows an EL-labeling of the poset from Figure 1.2-(a). Any linear ordering of the facets of $\Delta(P)$ compatible with the lexicographic order of the label sequences of the corresponding maximal chains of $P$ will be a shelling order for $\Delta(P)$. The shelling order of the facets of $\Delta(P)$ shown in Figure 1.6-(b) corresponds to the lexicographic ordering $(1, 2) < (2, 1) < (3, 1)$ of the label sequences for the maximal chains shown in Figure 1.6-(a).

<div align="center">(a)</div>

<div align="center">(b)</div>

Figure 1.6: An EL-labeling and the associated shelling of the order complex.

So, if there exists an EL-labeling of $\mathcal{E}(P)$, then $\Delta(P)$ is shellable and we say that $P$ is *EL-shellable*.

**Example 1.3.6.** The lattice of set partitions $\Pi_n$ is EL-shellable. The EL-labeling we give here first appears in [11] and is due to Ira Gessel. If $x \lessdot y$ in $\Pi_n$, then $y$ is obtained from $x$ by merging two blocks $B_1$ and $B_2$. Then the labeling

$$\lambda(x, y) = \max\{\min B_1, \min B_2\}$$

is an EL-labeling of $\mathcal{E}(\Pi_n)$. For example, the unique increasing chain in the interval $[\hat{0}, \hat{1}]$ in $\Pi_4$ is

$$\hat{0} = 1|2|3|4 \lessdot 12|3|4 \lessdot 123|4 \lessdot 1234 = \hat{1},$$

and has label sequence

$$\lambda(1|2|3|4, 12|3|4), \lambda(12|3|4, 123|4), \lambda(123|4, 1234) = 2, 3, 4.$$

### 1.3.2 Shellabilty and Homology Calculations

We can calculate the Betti numbers (Definition 1.1.10) $\beta_k(\Delta)$ of a simplicial complex $\Delta$ by understanding how a shelling order puts $\Delta$ together in a fashion that lets us explicitly understand

<div align="center">17</div>

the topology of $\Delta$. An $r$-dimensional facet $F_k$ of $\Delta$ is a *homology $r$-facet* if $F_k$ satisfies

$$\partial F_k = F_k \cap \bigcup_{i<k} \overline{F_i}.$$

So $F_k$ is a homology $r$-facet when $F_k$ attaches to $\Delta$ along its whole boundary in a shelling order. The Betti numbers of any shellable simplicial complex have a natural interpretation in terms of homology facets: the number of $r$-spheres in the homotopy type of $\Delta$ is the number of homology $r$-facets, as described in Theorem 1.3.3. In other words, $\beta_r(\Delta)$ is equal to the number of $r$-spheres in the homotopy type of $\Delta$.

**Example 1.3.7.** Recall from Example 1.3.4 that the non-pure simplicial complex $\Delta_b$ in Figure 1.1-(b) is shellable with shelling order $F_1 = \{1,2,4\}$, $F_2 = \{2,3\}$, $F_3 = \{2,5\}$, and $F_4 = \{3,5\}$. It is clear that $\Delta_b$ has the homotopy type of a 1-sphere. The boundary of the facet $F_4$, which is the pair of vertices $\{3\}$ and $\{5\}$, is contained in the union $\bigcup_{i<4} \overline{F_i}$, but no other facets have boundaries contained in the union of earlier facets.

## 1.4  Geometry of Distance-Based Phylogenetic Methods

A *phylogeny* is a mathematical model of the common evolutionary history of a group of genes or species. Phylogenies are often represented using a phylogenetic tree, defined formally in Definition 1.4.1. A phylogenetic tree modeling the evolutionary history of four species labeled "Dog, Cat, Frog, and Fish" is shown in Figure 1.7.

Phylogenetic trees are inferred from biological data such as DNA sequences, morphology, and fossil evidence. This thesis discusses distance-based methods that utilize numerical inputs, and our mathematical analysis does not depend on the methods by which these inputs are obtained. For context, it is worth noting that inputs for distance-based methods are often computed from amino acid or DNA sequence data using a statistical model of sequence evolution. There are many types of phylogenetic inference that use many types of data. We direct the interested reader to *Phylogenetics* by Charles Semple and Mike Steel [57] for a mathematical introduction to the field of phylogenetics and to *Inferring Phylogenies* by Joe Felsenstein [36] for an introduction to phylogenetics with a more thorough explanation of the biological perspective.

We define a *graph* as an ordered pair $G = (V(G), E(G))$ where $V(G)$ is the set of *vertices* of $G$, and $E(G) \subset \{(u,v) \mid u \neq v, \{u,v\} \subset V\}$ is the set of *edges* of $G$. Note that this definition precludes the existence of loops. Furthermore, we assume for the purposes of this thesis that $G$ has no multiple edges. A *path* $P \subset V(G)$ is a set of vertices $\{v_0, \ldots, v_n\}$ where $\{v_{i-1}, v_i\} \in E(G)$ for all $i \in [n]$. A graph $G$ is *connected* if there exists a path between every

Figure 1.7: A phylogeny for four species.

pair of vertices $\{u, v\} \subset V(G)$. A path $\Gamma$ in $G$ is a *cycle* if $v_0 = v_n$ and no vertex $v \in V(G)$ appears more than once in the sequence in $\Gamma$ except for $v_0$. A graph is *acyclic* if it contains no cycles. A *tree* is a connected, acyclic graph. The *degree* of a vertex $w$, denoted $\deg(w)$, is the number $|\{(u, v) \in E(G) \mid w \in \{u, v\}\}|$. The degree-one vertices of a tree graph are called *leaves*.

**Definition 1.4.1.** A *phylogenetic $X$- tree $T$* on the label set $X$ is an ordered pair $T = (\mathcal{T}, \phi)$ where $\mathcal{T}$ is a tree graph $\mathcal{T} = (E(\mathcal{T}), V(\mathcal{T}))$ with no vertices of degree two, and $\phi : X \to L(\mathcal{T})$ is a bijection from the label set $X$ onto the leaves $L(\mathcal{T})$ of $\mathcal{T}$. Two phylogenetic trees $T_1 = (\mathcal{T}_1, \phi_1)$ and $T_2 = (\mathcal{T}_2, \phi_2)$ are *isomorphic* if there exists a bijection $\varphi : V(\mathcal{T}_1) \to V(\mathcal{T}_2)$ that induces a bijection $E(\mathcal{T}_1) \to E(\mathcal{T}_2)$ and satisfies $\phi_2 = \varphi \circ \phi_1$.

It is often convenient to identify the set of species or genes $X$ that is being modeled as $[n] = \{1, 2, \ldots, n\}$ where $n = |X|$. The set $X$ may be referred to as *taxa* or *taxonomic units*. Two isomorphic phylogenetic $X$-trees with $X = [5]$ are shown in Figure 1.8. For the purposes of this thesis, we regard isomorphic phylogenetic $X$-trees as being equivalent, and omit the mention of the leaf or taxon label set $X$ when it is obvious from the context. Furthermore, it is often sufficient to write $T$ to denote a phylogenetic tree instead of using the ordered pair notation $T = (\mathcal{T}, \phi)$. A phylogenetic $X$-tree $T = \{\mathcal{T}, \phi\}$ is *rooted* if there exists a distinguished vertex $\rho$ of degree at least two such that $\rho \notin \phi(X)$. The two trees in Figure 1.8 are unrooted. The trees in Figure 1.7 and 1.9 are rooted.

Figure 1.8: Two isomorphic phylogenetic [5]-trees.

**Definition 1.4.2.** Two phylogenetic $X$-trees $T_1$ and $T_2$ have the same *tree topology* if their underlying trees $\mathcal{T}_1$ and $\mathcal{T}_2$, ignoring leaf labels, are isomorphic. We also say $T_1$ and $T_2$ have the same *shape* or *tree shape* in this case.

We can view a rooted tree as a directed graph where each edge $e \in E(T)$ is directed away from the root. In a rooted tree, we say a vertex $v \in V(T)$ is a *descendant* or *child* of $u \in V(T)$ if the path from $\rho$ to $v$ includes $u$. This digraph interpretation of a rooted tree induces a partial order $\leq_T$ on the vertices of the tree $T$ so that $u \leq_T v$ when $u$ is on the path from the root $\rho$ to $v$. Accordingly, if $u \leq_T v$, then $v$ is a descendant of $u$.

Let $\mathring{V}(T)$ denote the interior (i.e. non-leaf) vertices of $T$. We say that $T = (\mathcal{T}, \phi)$ is *binary* if every $u \in \mathring{V}(T)$ has degree 3, except the root $\rho$, which if it exists, has degree 2. The motivation for the term "binary" in this definition is that every interior vertex, including $\rho$ if $T$ has a root, has exactly 2 descendants. However, the term binary may be applied to unrooted trees, in which every interior vertex has degree 3. We associate a special term to interior vertices that either have degree greater than or equal to 4 or have more than two descendants.

**Definition 1.4.3.** A *polytomy* is a vertex in a phylogenetic tree that is either a root $\rho$ of a rooted tree satisfying $\deg \rho \geq 3$ or an interior vertex $u$ that is not a root satisfying $\deg u \geq 4$.

A phylogenetic tree with a polytomy vertex models an evolutionary history where there is either a multi-way speciation event or where it is understood that there is insufficient data to

completely resolve the history. The two trees in Figure 1.8 each have a polytomy vertex.

In a rooted tree, there may be more than one ordering on the internal vertices that may correspond to speciation events in the inferred evolutionary history. To make this precise, we make the following definition:

**Definition 1.4.4.** A *rank function* on $T$ is a bijection $r : \mathring{V}(T) \to \{1, 2, \ldots, |\mathring{V}(T)|\}$ satisfying $u \leq_T v \Rightarrow r(v) \leq r(u)$.

The number of rank functions on $T$ is given by the next formula, which appears in [66] as Exercise 1 in the supplementary exercises for Chapter 3.

$$|\mathring{V}(T)|! / \prod_{v \in V(\mathring{T})} |\mathrm{de}(v)| \tag{1.2}$$

where $\mathrm{de}(v)$ denotes the set of descendants of $v$.

A rooted phylogenetic tree with a rank function is called a *ranked phylogenetic tree*. Sometimes it is useful to consider a specific type of subgraph of the tree $T$:

**Definition 1.4.5.** Let $T$ be a rooted phylogenetic tree. Let $v$ be a vertex in $V(T)$. Then the *clade* associated to $v$ is the sub-tree of $T$ consisting of $v$ along with the set of all descendants of $v$.

**Example 1.4.6.** For the rooted, ranked phylogenetic [5]-tree $T = \{\mathcal{T}, \phi : [5] \to L(\mathcal{T})\}$ in Figure 1.9, let $u_i$ denote the interior vertex of rank $i$, i.e. $r(u_i) = i$. Then the clade associated to $u_3$ is the subtree formed by the induced subgraph with vertex set $\{u_1, u_2, u_3, \phi(1), \phi(2), \phi(3), \phi(4)\}$.



Figure 1.9: A ranked phylogenetic [5]-tree.

Note that ranked phylogenetic [n]-trees are naturally in bijection with maximal chains in the lattice of set partitions $\Pi_n$, which was defined in Example 1.1.4.

**Example 1.4.7.** The maximal chain

$$\pi_5 = 1|2|3|4|5 \lessdot 12|3|4|5 \lessdot 12|34|5 \lessdot 1234|5 \lessdot 12345 = \pi_1$$

in $\Pi_5$ corresponds to the ranked phylogenetic tree in Figure 1.9.

### 1.4.1 Distance-Based Phylogenetic Methods

Distance-based phylogenetic reconstruction methods get their name from the characteristics of their inputs and outputs, which are types of pairwise distance functions on the set of species being modeled. In particular, all distance-based methods take dissimilarity maps as inputs and return tree metrics as outputs. We formally define both terms now:

**Definition 1.4.8.** A function $\delta : X \times X \to \mathbb{R}$ is a *dissimilarity map* on $X$ if

1. $\delta(x, y) = \delta(y, x)$ for all $\{x, y\} \subset X$, and

2. $\delta(x, x) = 0$ for all $x \in X$.

If we identify $X$ as $[n]$, and label the coordinates of $\mathbb{R}^{\binom{n}{2}}$ with the two-element subsets of $[n]$, every point in $\mathbb{R}^{\binom{n}{2}}$ can be identified as a dissimilarity map on $[n]$.

**Definition 1.4.9.** A function $d : X \times X \to \mathbb{R}$ is a *tree metric* if there exists a phylogenetic $X$-tree $T$ and a function $w : E(T) \to \mathbb{R}_{\geq 0}$ such that $d(x, y) = \sum_{e \in P(x,y)} w(e)$ where $P(x, y)$ is the unique path in $\mathcal{T}$ between the leaves $\phi(x)$ and $\phi(y)$. We refer to the pair $(T, w)$ as a *weighted tree.*

We say that $d$ is *realized* by the phylogenetic tree $T$. Note that $d$ may be realized by a rooted tree such as the tree in Figure 1.9, or an unrooted tree such as the trees in Figure 1.8. There are interesting characterizations of tree metrics in terms of constraints on the distances between the leaves of the tree, which we now describe.

**Definition 1.4.10.** A dissimilarity map $\delta$ on $X$ satisfies the *four-point condition* if, for every four (not necessarily distinct) elements $\{w, x, y, z\} \subset X$,

$$\delta(w, x) + \delta(y, z) \leq \max\{\delta(w, y) + \delta(x, z), \delta(w, z) + \delta(x, y)\}$$

The following theorem is due independently to K. A. Zaretskii [71], J. M. S. Simões-Pereira [59], and Peter Buneman [18], [19].

**Theorem 1.4.11.** *A dissimilarity map $\delta$ on $X$ satisfies the four-point condition if and only if $\delta$ is a tree metric on $X$.*

If we take $w = x$ we can recover the traditional triangle inequality from the four-point condition on $X$. Therefore the set of tree metrics on $X$ is a subset of the set of all metrics on $X$.

We now introduce a subset of the set of tree metrics that are commonly associated to tree metrics realized by rooted trees. For a phylogenetic $X$-tree $T = (\mathcal{T}, \phi)$, we usually discuss tree metrics as a function on the set $X$ only, but we can extend the function $d$ to all vertices of $\mathcal{T}$ by taking $d(u, v)$ to be the sum of edge weights on the unique path from $u$ to $v$ for any pair of vertices $\{u, v\} \subset V(\mathcal{T})$. We do so to make Definition 1.4.12:

**Definition 1.4.12.** An edge-weighting of a rooted tree $T$ with root vertex $\rho$ is *equidistant* if the tree metric $d$ realized by the weighted tree $(T, w)$ satisfies

1. $d(\rho, \phi(x)) = d(\rho, \phi(y))$ for all $x, y \in X$, and

2. $d(v, \phi(x)) \leq d(u, \phi(x))$ whenever $u \leq_{\mathcal{T}} v \leq_{\mathcal{T}} \phi(x)$.

Tree metrics realized by rooted trees with equidistant weightings also have an alternate characterization. To make this characterization, we first give Definition 1.4.13:

**Definition 1.4.13.** A dissimilarity map $\delta : X \times X \to \mathbb{R}$ is an *ultrametric* on $X$ if for every three distinct elements $\{x, y, z\} \subset X$,

$$\delta(x, y) \leq \max\{\delta(x, z), \delta(y, z)\}$$

The following theorem appears as part of Theorem 7.2.5 in *Phylogenetics* by Charles Semple and Mike Steel [57]:

**Theorem 1.4.14.** *Let $\delta$ be a dissimilarity map on $X$. Then $\delta$ is an ultrametric if and only if there exists a rooted phylogenetic $X$-tree $T$ with an equidistant weighting $w$ such that the weighted tree $(T, w)$ realizes $\delta$.*

As we mentioned after Definition 1.4.8, we can identify the set of all dissimilarity maps on $[n]$ as $\mathbb{R}^{\binom{n}{2}} = \mathbb{R}^{n(n-1)/2}$. The definitions of tree metrics and ultrametrics are less general, but since tree metrics and ultrametrics on the set $[n]$ are also dissimilarity maps, they can be identified as points in the same space $\mathbb{R}^{n(n-1)/2}$. Because of Theorems 1.4.11 and 1.4.14, we can think of these points either as points corresponding to certain weighted trees or points in certain polyhedral subsets of $\mathbb{R}^{n(n-1)/2}$:

**Definition 1.4.15.** Let $\mathcal{T}_n$ denote the set of all tree metrics on $[n]$ and let $\mathcal{ET}_n$ denote the set of all ultrametrics on $[n]$. Then $\mathcal{T}_n$ and $\mathcal{ET}_n$ are proper subsets of $\mathbb{R}^{n(n-1)/2}$ that we identify as *tree space* and *equidistant tree space*, respectively.

Note that $\mathcal{T}_n$ and $\mathcal{ET}_n$ provide a different notion of tree space than the tree space described in the paper "Geometry of the space of phylogenetic trees" by Louis Billera, Susan Holmes, and Karen Vogtmann [8]. The spaces $\mathcal{T}_n$ and $\mathcal{ET}_n$ are well-known to be polyhedral fans (Definition 1.2.8) in $\mathbb{R}^{n(n-1)/2}$ of dimension $2n-3$ [62] and $n-1$ [3] respectively.

**Definition 1.4.16.** A *distance-based phylogenetic method* or simply a *distance-based method* is any function $f$ that takes a dissimilarity map $\delta$ as an input and returns a tree metric $d$, where $d$ is either an arbitrary tree metric in $\mathcal{T}_n$ or an ultrametric in $\mathcal{ET}_n$, as an output.

A distance-based phylogenetic method $f$ induces a partition of the space $\mathbb{R}^{n(n-1)/2}$ of all dissimilarity maps into a family of regions

$$\{C(T) : f(x) \text{ is a tree metric realized by the combinatorial tree } T \text{ for all } x \in C(T)\}.$$

Each region $C(T)$ is the set of all dissimilarity maps that the distance-based method sends to a cone in the polyhedral fan corresponding to $T$. If we only consider inputs with positive entries, we can consider a partition of the *positive orthant* $\mathbb{R}_{\geq 0}^{n(n-1)/2}$ of $\mathbb{R}^{n(n-1)/2}$. In Chapter 4, we study the partition of $\mathbb{R}_{\geq 0}^{n(n-1)/2}$ induced by the UPGMA algorithm (Algorithm 4.2.1). We will see that this partition is a family of polyhedral cones whose $\mathcal{V}$-representation is described in terms of maximal chains in the lattice of set partitions $\Pi_n$ introduced in Example 1.1.4. In Chapter 5, we use polyhedral geometry to investigate regions in the partition of the input space $\mathbb{R}^{n(n-1)/2}$ by distance-based phylogenetic methods corresponding to families of binary resolutions of tree metrics realized by trees with polytomies.

# Chapter 2

# Lexicographic Shellability and the Characterization of Posets

## 2.1 Introduction

In this chapter we discuss the notion of lexicographic shellability and investigate the problem of characterizing posets by the types of lexicographic shellings they admit. The concept of lexicographic shellability was first introduced by Anders Björner in [11]. We explained the approach of using one type of lexicographic shelling known as an EL-labeling (Definition 1.3.5) to shell the order complex of a poset in Chapter 1, and we will explain another type of lexicographic shellability induced by CL-labelings (Definition 2.3.1) in this chapter. There are many examples of lexicographically shellable posets in the literature, often with a view towards showing that fundamental types of posets admit lexicographic shellings. Many lexicographically shellable posets of interest in the literature, including those studied in this chapter, are lattices (Definition 2.2.1).

For example, Björner showed that geometric lattices were lexicographically shellable in [11], Adriano Garsia showed that upper semimodular lattices are lexicographically shellable in [40], Richard Stanley gave a lexicographic shelling of supersolvable lattices in [64], and Michelle Wachs and James Walker showed that geometric semilattices are lexicographically shellable in [69]. John Shareshian showed that the subgroup lattice of a finite solvable group is CL-shellable (see Definition 2.3.1) in [58], and Russ Woodroofe showed that this lattice is EL-shellable in [70]. We briefly explore the relationship between EL-shellability and CL-shellability of finite graded atomic lattices in Section 2.3.

In Section 2.4 we characterize geometric lattices (Definition 2.2.4) as those finite, atomic lattices such that every atom ordering induces a lexicographic shelling of the order complex given by an edge-labeling of the Hasse diagram known as a minimal labeling (Definition 2.2.3). Equiv-

alently, we show geometric lattices are exactly those finite lattices such that every ordering on the join-irreducibles (Definition 2.2.2) induces a lexicographic shelling. This new characterization fits into a similar paradigm as Peter McNamara's characterization of supersolvable lattices as those lattices admitting a different type of lexicographic shelling, namely one in which each maximal chain is labeled with a permutation of $\{1, \ldots, n\}$ [49]. In Section 2.5 we give a similar characterization of semimodular lattices. The results in Sections 2.4 and 2.5 are joint work with Patricia Hersh and appear in the paper [23].

## 2.2  Lattices and Minimal Labelings

Lattices are of particular interest in the theory of posets. Every finite lattice admits a type of edge-labeling known as a minimal labeling. We now briefly develop some terminology related to lattices, as well as some additional terminology for posets necessary to define minimal labelings.

If $p$ and $q$ are elements of a poset $P$ and there exists $r \in P$ such that $r \geq p$ and $r \geq q$, then $r$ is an *upper bound* for $p$ and $q$. We say $p$ and $q$ have a *least upper bound* or *join* if there exists an upper bound $s$ of $p$ and $q$ such that $s \leq r$ for all other upper bounds $r$ of $p$ and $q$. We write $s = p \vee q$ if $s$ is the least upper bound of $p$ and $q$ in $P$. Similarly, we say $s$ is the *greatest lower bound* or *meet* of $p$ and $q$ in $P$ and write $s = p \wedge q$ if $s \geq r$ for all lower bounds $r$ of $p$ and $q$ in $P$.

**Definition 2.2.1.** A poset $L$ is a *lattice* if every pair of elements $p$ and $q$ in $L$ have both a join and a meet.

Consequently, every finite lattice must have a $\hat{0}$ and a $\hat{1}$.

**Definition 2.2.2.** An element $r$ of a poset $P$ is *join-irreducible* if we cannot write $r = p \vee q$ for any $p$ and $q$ in $P$ satisfying $p \neq q$, $p < r$, and $q < r$.

An *atom* is an element of $P$ that covers $\hat{0}$. In a finite poset, the join-irreducibles are precisely the elements that cover exactly one element, and include the atoms. For a poset (or lattice) $P$ we write $\mathrm{JoinIrred}(P)$ to denote the set of join-irreducibles of $P$, and $A(P)$ to denote the set of atoms of $P$. A lattice $L$ is *atomic* is every element of $L$ is the join of atoms. Note that if $L$ is an atomic lattice, $\mathrm{JoinIrred}(L) = A(L)$.

**Definition 2.2.3.** Given any bijection $\gamma : \mathrm{JoinIrred}(L) \to [n]$, the map $\gamma$ induces a *minimal labeling* $\lambda_\gamma : E(L) \to [n]$ by the rule

$$\lambda_\gamma(x, y) = \min\{\gamma(j) : j \in \mathrm{JoinIrred}(y) \setminus \mathrm{JoinIrred}(x)\}.$$

A finite poset $P$ is *graded* if every maximal chain has the same number of elements in it; in this case, there is a rank function $\rho$ defined recursively by $\rho(x) = 0$ for $x$ a minimal element of $P$ and $\rho(y) = \rho(x) + 1$ for $x \lessdot y$. A graded lattice is *semimodular* if it has a rank function $\rho$ that satisfies

$$\rho(x \wedge y) + \rho(x \vee y) \leq \rho(x) + \rho(y). \tag{2.1}$$

**Definition 2.2.4.** A finite lattice is *geometric* if it is atomic and semimodular.

Proposition 2.2.5 appears as Corollary 1, p. 81 in Garrett Birkhoff's book *Lattice Theory* [10], and gives an alternate formulation of semimodularity for graded lattices that will be convenient in proofs that appear later in this chapter.

**Proposition 2.2.5** (Birkhoff). *Let $L$ be a finite lattice. The following two conditions are equivalent:*

1. *$L$ is graded, and the rank function $\rho$ of $L$ satisfies the semimodularity condition (2.1) above.*

2. *If $x$ and $y$ both cover $x \wedge y$, then $x \vee y$ covers both $x$ and $y$.*

Geometric lattices are an important class of lattices. For example, the class of geometric lattices includes all of the intersection lattices of real, central hyperplane arrangements. Also, every geometric lattice is the lattice of flats of a *matroid*, which we define next. Note that matroids have many equivalent characterizations. The one we provide, Definition 2.2.6, appears in Richard Stanley's article [63] as part of a discussion on hyperplane arrangements.

**Definition 2.2.6.** Let $S$ be a finite set and let $2^S = \{T : T \subseteq S\}$. A (finite) *matroid* is a pair $M = (S, \mathcal{J})$, where $\mathcal{J}$ is a collection of subsets of $S$ satisfying the following two conditions:

- $\mathcal{J}$ is a nonempty abstract simplicial complex (Definition 1.1.1), i.e. $\mathcal{J} \neq \emptyset$ and if $J \in \mathcal{J}$ and $I \subseteq J$, then $I \in \mathcal{J}$.

- For any $T \subset S$, the inclusion-wise maximal elements of $\mathcal{J} \cap 2^T$ have the same cardinality.

The collection $\mathcal{J}$ is called the collection of *independent sets*. An independent set $T$ is a *basis* if $T \cup \{s\} \notin \mathcal{J}$ for any $s \in S$. The *rank* of a set $T \subseteq S$ is the number

$$\text{rank}(T) = \max\{|I| : I \in \mathcal{J} \text{ and } I \subseteq T\}.$$

So, $\text{rank}(\emptyset) = 0$ and $\text{rank}(M) = \text{rank}(S)$. A *$k$-flat* is a maximal subset of rank $k$. The *closure* $\overline{T}$ is the smallest flat containing $T$, or in other words,

$$\overline{T} = \bigcap_{\text{flats } F \supseteq T} F.$$

We can identify geometric lattices with matroids in the following manner: if $S$ is the set of atoms of a geometric lattice $L$, then $T \subseteq S$ is independent if and only if the join of $T$ in $L$ has rank equal to the cardinality of $T$. The next example relates the concepts of matroid and geometric lattice via a lattice we first saw in Chapter 1.

**Example 2.2.7.** The lattice of set partitions $\Pi_n$ from Example 1.1.4 is a geometric lattice. One can see this by checking that $\Pi_n$ is finite, atomic, and semimodular. The atoms of $\Pi_n$ are the partitions consisting of exactly one block containing two elements and all other blocks consisting of singleton elements. For example, the atoms of $\Pi_4$ are the partitions

$$12|3|4, 13|2|4, 14|2|3, 23|1|4, 24|1|3, 34|1|2.$$

The *graphical matroid* of a finite graph $G$ is the matroid $M(G) = (S, \mathcal{J})$ where $S$ is the set $E(G)$ of edges of $G$ and $\mathcal{J}$ is the collection of subsets of $E(G)$ that are *forests* of $G$, or collections of connected acyclic subgraphs. Let $\pi$ be a set partition of the vertices of $G$ and let $F_\pi$ be the subset of $E(G)$

$$F_\pi = \{e \in E(G) : \text{both endpoints of } e \text{ are contained in the same block of } \pi\}.$$

Then $F_\pi$ is a flat of $M(G)$. Consider the case $G = K_n$, and note that every set partition of $[n]$ can be obtained by a subset of $E(K_n)$. So, $\Pi_n$ is isomorphic to the lattice of flats of the graphical matroid of the graph $K_n$, the complete graph on $n$ vertices.

The lattice $\Pi_n$ can also be viewed as the intersection lattice of the real hyperplane arrangement consisting of hyperplanes of the form $x_i = x_j$, $\{i, j\} \in \binom{[n]}{2}$. This hyperplane arrangement is also known as the *Type A Coxeter arrangement*.

Since geometric lattices are atomic, we can induce a minimal labeling $\lambda_\gamma$ on a geometric lattice $L$ from any total ordering $\gamma$ of the set $A(L)$ of the atoms of $L$. Such minimal labelings for geometric lattices were originally introduced by Björner in [11]; the following theorem shows the motivation for introducing them.

**Theorem 2.2.8** (Björner [11]). *The minimal labeling resulting from any linear ordering of the atoms in a geometric lattice is an EL-labeling.*

In Section 2.4 we show that a lattice $L$ is geometric if and only if $\lambda_\gamma$ is an EL-labeling for every total ordering $\gamma$ of $A(L)$.

## 2.3 CL-Shellability and EL-Shellability

A concept closely related to EL-shellability is that of CL-shellability. The idea of a CL-shellable poset was introduced by Anders Björner and Michelle Wachs in [13] in order to show that

the Bruhat order poset for a Coxeter group is lexicographically shellable. To understand the notion of CL-shellability we must first introduce the idea of a chain-edge labeling: a labeling $\lambda : \mathcal{E}(P) \to \Lambda$, where $\Lambda$ is a poset, is a *chain-edge labeling* if whenever two chains $C, C'$ in $P$ coincide along there first $d$ edges, then their labels coincide along their first $d$ edges as well. In other words, whenever $x_i \lessdot x_{i+1}$ for $x_i, x_{i+1} \in C \cap C'$, $\lambda_C(x_i, x_{i+1}) = \lambda_{C'}(x_i, x_{i+1})$ for all $i \leq d - 1$. When a poset is given a chain-edge labeling, it is possible that an edge in $\mathcal{E}(P)$ may have more than one label depending on which chain we identify it as a member of.

Since it is possible that in a chain-edge labeling an edge in $\mathcal{E}(P)$ may have more than one label depending on which chain we identify it as a member of, we remove this ambiguity when working with chain-edge labelings by introducing the concept of rooting an interval in a poset. Given a poset $P$, if $[x, y] \subset P$ is an interval in $P$ and $r$ is a maximal chain from $\hat{0}$ to $x$, then the pair $([x, y], r)$ is a *rooted interval* with root $r$, denoted $[x, y]_r$. Note that this implies if $C$ is a maximal chain of $[x, y]$, then $C \cup r$ is a maximal chain from $\hat{0}$ to $y$.

**Definition 2.3.1.** A chain-edge labeling $\lambda$ of the edges $\mathcal{E}(P)$ of the Hasse diagram of a poset $P$ is called a *CL-labeling* (chain lexicographic labeling) if for every rooted interval $[x, y]_r$ in $P$,

1. there is a unique increasing maximal chain $C$ in $[x, y]_r$, and

2. the label sequence of $C$ is lexicographically smaller than the label sequence of every other saturated chain in the interval $[x, y]_r$.

If $P$ admits a CL-labeling, we say $P$ is *CL-shellable*. Since every EL-labeling is a CL-labeling, EL-shellability implies CL-shellability. It is unknown, in general, if CL-shellability implies EL-shellability. Section 4.1 of the article [68] by Michelle Wachs contains a thorough discussion of the relationships between different types of lexicographic shellability for posets.

### 2.3.1 Recursive Atom Orderings

The notion of a recursive atom ordering was introduced by Björner and Wachs in [14]. Recursive atom orderings give an alternative formulation of CL-shellability for finite posets that does not depend on an edge-labeling. Recall that the *length* of a finite poset $P$ is the maximum length of a chain of $P$.

**Definition 2.3.2.** A finite poset $P$ is said to admit a *recursive atom ordering* if either its length $\ell(P) = 1$ or $\ell(P) > 1$ and there is an ordering $a_1, a_2, \ldots, a_t$ of the atoms of $P$ that satisfies the following two conditions:

1. For all $j \in [t]$, the interval $[a_j, \hat{1}]$ admits a recursive atom ordering in which the atoms of $[a_j, \hat{1}]$ that belong to $[a_i, \hat{1}]$ for some $i < j$ come first.

2. For all $i < j$, if $a_i, a_j < y$ then there is a $k < j$ and an atom $z$ of $[a_j, \hat{1}]$ such that $a_k < z \le y$.

Figure 2.1 shows a recursive atom ordering of a poset. The atoms of the poset are ordered left to right and labeled $A1, A2, A3$, and $A4$. The rank 2 elements are ordered left to right, and labeled $B1, B2$, and $B3$. The ordering of the rank two elements induces their ordering as atoms of the intervals $[Ai, \hat{1}]$.



Figure 2.1: A recursive atom ordering.

The following theorem appears in [14] as Theorem 3.2:

**Theorem 2.3.3** (Björner and Wachs). *A graded poset $P$ admits a recursive atom ordering if and only if $P$ is CL-shellable.*

A consequence of Theorem 2.3.3 is that if a finite graded atomic lattice $L$ has a recursive atom ordering that induces a minimal labeling that is an EL-labeling, then $L$ is both CL-shellable and EL-shellable. This suggests the following conjecture:

**Conjecture 2.3.4.** *Let $L$ be a finite graded atomic lattice. Then if $L$ is CL-shellable, $L$ is EL-shellable.*

Conjecture 2.3.4 posits the equivalence of EL- and CL-shellability for finite graded atomic lattices. One approach to proving this conjecture would be as follows: given a finite graded

30

atomic lattice that is CL-shellable, we know that $L$ admits a recursive atom ordering. Since $L$ is an atomic lattice, this ordering induces a minimal labeling. If one could show this labeling was always an EL-labeling, this would establish Conjecture 2.3.4. Unfortunately, the minimal labeling induced by the recursive atom ordering of the lattice in Figure 2.1 is not an EL-labeling, as shown in Figure 2.2.



Figure 2.2: A minimal labeling induced by a recursive atom ordering.

While this does not show that Conjecture 2.3.4 is false, it does show that this method of proof will not work.

## 2.4 Lexicographic Shellability Characterizations of Geometric Lattices

Peter McNamara proved that supersolvable lattices can be characterized as lattices admitting a certain type of EL-labeling known as an $S_n$-EL-labeling [49]. Each maximal chain is labeled by the set of labels $\{1, \ldots, n\}$ with each label used exactly once in each maximal chain. Previously, Richard Stanley had proven that all supersolvable lattices admit such EL-labelings in [64]. Thus, McNamara's result gave a new characterization of supersolvable lattices: that a finite lattice is supersolvable if and only if it has an $S_n$-EL-labeling.

This section is devoted to giving two new characterizations of geometric lattices. The first is based on atom orderings for finite atomic lattices. The second characterization replaces this with a condition on the orderings of the join-irreducibles for finite lattices so as to avoid assuming a priori that the lattices are atomic. In both cases, we prove for any finite lattice $L$ that if every ordering of the join-irreducibles induces a minimal labeling which is an EL-labeling, then $L$ is a geometric lattice. To this end, we first develop some helpful properties of minimal labelings.

**Lemma 2.4.1.** *Let $L$ be a finite atomic lattice and let $\lambda_\gamma$ be a minimal labeling on $L$. Then $x_i \lessdot x_{i+1} \leq x_j \lessdot x_{j+1}$ in $L$ implies $\lambda_\gamma(x_i, x_{i+1}) \neq \lambda_\gamma(x_j, x_{j+1})$. In other words, the labels on any particular saturated chain are distinct.*

*Proof.* This is immediate from the fact that $A(x_{j+1}) \setminus A(x_j)$ is disjoint from $A(x_{i+1}) \setminus A(x_i)$ for $i \neq j$. $\qquad \square$

**Lemma 2.4.2.** *Let $L$ be a finite lattice. Then $\mathrm{JoinIrred}(u) \subseteq \mathrm{JoinIrred}(v)$ if and only if $u \leq v$. Moreover, $u = v$ if and only if $\mathrm{JoinIrred}(u) = \mathrm{JoinIrred}(v)$. In the special case of a finite atomic lattice $L$ we have $A(u) \subseteq A(v)$ if and only if $u \leq v$, and we have $A(u) = A(v)$ if and only if $u = v$.*

*Proof.* This follows from two facts: (1) that every element of a finite lattice $L$ is a join of join-irreducibles, and (2) that the only join-irreducibles in an atomic lattice are the atoms. $\qquad \square$

**Lemma 2.4.3.** *Let $L$ be a finite lattice and suppose that there exist $x, y \in L$ that both cover $x \wedge y$, but with $x$ not covered by $x \vee y$. Then for $j$ any join-irreducible satisfying $y = (x \wedge y) \vee j$, we have that $j \notin \mathrm{JoinIrred}(z)$ for any $z$ covering $x$.*

*Proof.* Assume by way of contradiction that the join-irreducible $j$ given above satisfies $j \leq z$ for some $x \lessdot z$. Note that $x \wedge y \leq x \lessdot z$, which together with $j \leq z$ implies $(x \wedge y) \vee j \leq z$. But $(x \wedge y) \vee j = y$, so we may conclude that $x \vee y \leq z$. This contradicts the fact that $x \vee y$ does not cover $x$, completing our proof. $\qquad \square$

Now to our first characterization of geometric lattices.

**Theorem 2.4.4.** *Let $L$ be a finite atomic lattice. Then $L$ is geometric if and only if every atom ordering induces a minimal labeling that is an EL-labeling.*

*Proof.* Björner proved in Theorem 2.2.8 that all of the minimal labelings for a geometric lattice are EL-labelings. We now prove the converse. Since we assume that $L$ is atomic, what remains is to prove that $L$ is semimodular.

Suppose otherwise. By Proposition 2.2.5, there must exist $x, y \in L$ such that $x$ and $y$ both cover $x \wedge y$ but $x \vee y$ does not cover $x$. By Lemma 2.4.2, we may choose some atom

32

$a_x \in A(x) \setminus A(x \wedge y)$ such that $a_x \notin A(y)$. Since $L$ is an atomic lattice, there must also exist $a_y \in A(y)$ such that $(x \wedge y) \vee a_y = y$. By Lemma 2.4.3, $a_y \notin A(z)$ for any $z$ such that $x \lessdot z$. This implies $a_y \neq a_x$, since $a_y \notin A(z)$ for all $z$ satisfying $x \lessdot z$, which in particular implies $a_y \notin A(x)$.

Now consider any atom ordering $\gamma : \mathcal{A}(L) \to [n]$ such that $\gamma(a_x) = 1$ and $\gamma(a_y) = 2$. Since $a_x \in A(x) \setminus A(x \wedge y)$ and $\gamma(a_x) = 1$, we know that $\lambda_\gamma(x \wedge y, x) = 1$. Let

$$C := \quad x \wedge y = x_0 \lessdot x = x_1 \lessdot x_2 \lessdot \cdots \lessdot x_k = x \vee y$$

be the lexicographically smallest saturated chain on the interval $[x \wedge y, x \vee y]$. By Lemma 2.4.3, $a_y \notin A(x_2)$. Therefore, $\lambda_\gamma(x_1, x_2) \neq 2$. By Lemma 2.4.1, there is no repetition in the label sequence, implying $\lambda_\gamma(x_1, x_2) > 2$. For some $2 < j \leq k$, we must have $a_y \in A(x_j) \setminus A(x_{j-1})$, implying $\lambda_\gamma(x_{j-1}, x_j) = 2$.

But $\min\{\gamma(a) | a \in A(x_2) \setminus A(x_1)\} \geq 3$, so $\lambda_\gamma(x_1, x_2) > \lambda_\gamma(x_{j-1}, x_j)$. This implies that $C$ cannot have weakly increasing labels, hence that $\lambda_\gamma$ is not an EL-labeling. $\qquad \square$

Next we give a closely related alternative characterization of geometric lattices which avoids making the assumption that the lattices are atomic. The essence of the proof will be a reduction to the atomic case. We thank Peter McNamara for pointing out that one does not need to assume that a finite lattice is atomic in order to give a lexicographic shellability characterization of geometric lattices.

**Theorem 2.4.5.** *A finite lattice $L$ is a geometric lattice if and only if every ordering of the join-irreducibles induces a minimal labeling $\lambda_\gamma$ which is an EL-labeling.*

*Proof.* One direction is well-known, so we focus on the other direction. That is, we will assume there is some join-irreducible that is not an atom, and use this to produce an ordering on join-irreducibles whose associated minimal labeling is not an EL-labeling. The case where all join-irreducibles are atoms has already been handled in Theorem 2.4.4.

Suppose there exists $v \in \mathrm{JoinIrred}(L)$ that is not an atom. In this case, we may choose such a $v$ so that if $a < v$ for $a \in \mathrm{JoinIrred}(L)$ then $a$ is an atom. It is well known (see [66], p. 286) that in a finite lattice, the join-irreducibles are precisely the elements that cover exactly one other element. Let $u$ be the unique element in $L$ with $u \lessdot v$. Thus, $\mathrm{JoinIrred}(u)$ is entirely composed of atoms and $|\mathrm{JoinIrred}(u)| = k$ for some $1 \leq k \leq n - 1$ where $n$ is the number of join-irreducibles in $L$. Consider an ordering $\gamma$ on the join-irreducibles such that $\gamma(v) = 1$ and $\{\gamma(x) | x \in \mathrm{JoinIrred}(u)\} = \{2, 3, \ldots, k + 1\}$. The lexicographically smallest label sequence for any saturated chain in the interval $[\hat{0}, v]$ must then have a descent, because all saturated chains must include $u$, but $\lambda_\gamma(u, v) = 1$ while $\lambda_\gamma(x, y) > 1$ for all covering relations $x \lessdot y$ in the interval $[\hat{0}, u]$. Thus, this minimal labeling $\lambda_\gamma$ is not an EL-labeling. $\qquad \square$

## 2.5 Lexicographic Shellability Characterization of Semimodular Lattices

Our next theorem, Theorem 2.5.1, characterizes semimodular lattices, and previously appeared in the paper [54] of Ivan Rival. We nonetheless include our proof of Theorem 2.5.1 both for its new approach to this sort of question and because it provides a significantly higher level of detail than appears in the argument in [54]. We thank Victor Reiner for suggesting the question of characterizing semimodular lattices using similar techniques to those we applied in Section 2.4. First we make an observation that the proof of Theorem 2.5.1 will rely upon.

*Observation.* Let $L$ be a finite lattice with $|\operatorname{JoinIrred}(L)| = n$ and let $x \in L$. If $|\operatorname{JoinIrred}(x)| = k$ and $\hat{\gamma} : \operatorname{JoinIrred}([\hat{0}, x]) \to [k]$ is a linear extension of the subposet of join-irreducibles of the interval $[\hat{0}, x]$, then there exists a linear extension $\gamma : \operatorname{JoinIrred}(L) \to [n]$ of the subposet $\operatorname{JoinIrred}(L)$ that restricts to the map $\hat{\gamma}$.

**Theorem 2.5.1.** *Let $L$ be a finite lattice with $|\operatorname{JoinIrred}(L)| = n$. Suppose that for every linear extension $\gamma : \operatorname{JoinIrred}(L) \to [n]$ of the subposet $\operatorname{JoinIrred}(L)$ of join-irreducibles in $L$, the resulting minimal labeling $\lambda_\gamma$ is an EL-labeling on $L$. Then $L$ is (upper) semimodular.*

*Proof.* Assume by way of contradiction that $x$ and $y$ cover $x \wedge y$ but that $x \vee y$ does not cover $x$. Lemma 2.4.2 shows that there exist join-irreducibles $j_x \in \operatorname{JoinIrred}(x) \setminus \operatorname{JoinIrred}(x \wedge y)$ and $j_y \in \operatorname{JoinIrred}(y) \setminus \operatorname{JoinIrred}(x \wedge y)$. Let $k$ be the number of elements in $\operatorname{JoinIrred}(x \wedge y)$. Notice that if $k = 0$, then $x$ and $y$ are atoms, so in particular $x$ and $y$ are join-irreducibles $j_x := x$ and $j_y := y$.

Now we choose a linear extension $\gamma$ of the subposet $\operatorname{JoinIrred}(L)$ of $L$ comprised of the join-irreducibles. By Observation 2.5, we may choose $\gamma$ so that it assigns exactly the values in $\{1, \ldots, k\}$ to the join-irreducibles in $[\hat{0}, x \wedge y]$. Moreover, we may insist that $\gamma(j_x) = k + 1$ and $\gamma(j_y) = k + 2$, choosing how $\gamma$ assigns the values in $\{k+3, \ldots, n\}$ to the subposet comprised of the remaining join-irreducibles by taking any linear extension of the remaining join-irreducibles.

Denote the lexicographically smallest maximal chain in the interval $[x \wedge y, x \vee y]$ by

$$C = x \wedge y \lessdot x \lessdot x_2 \lessdot \cdots x_{m-1} \lessdot x_m = x \vee y.$$

We have assumed that $x \vee y$ does not cover $x$, implying $m > 2$. Our constraints given above on our choice of $\gamma$ imply that $\lambda_\gamma(x_1, x_2) \notin \{1, \ldots, k+2\}$, since Lemma 2.4.3 ensures that $j_y \notin \operatorname{JoinIrred}(x_2)$. Thus, $\lambda_\gamma(x, x_2) \geq k + 3$. But since $\operatorname{JoinIrred}(y) \subset \operatorname{JoinIrred}(x \vee y)$, we then must have $j_y \in \operatorname{JoinIrred}(x_\ell) \setminus \operatorname{JoinIrred}(x_{\ell-1})$ for some $2 < \ell \leq m$. This implies $\min(\{\gamma(j)|j \in \operatorname{JoinIrred}(x_\ell) \setminus \operatorname{JoinIrred}(x_{\ell-1})\}) \leq k + 2$. Hence, $\lambda_\gamma(x_{\ell-1}, x_\ell) \leq k + 2 < r = \lambda_\gamma(x, x_2)$, forcing the chain $C$ to have a descent, a contradiction to this being an EL-labeling. Thus, knowing that

$x$ and $y$ both cover $x \wedge y$ does imply in our setting that $x \vee y$ covers both $x$ and $y$. Thus, $L$ is semimodular.

$\square$

## 2.6   Discussion

A concept closely related to that of an EL-labeling is that of an R-labeling. We give the definition from Chapter 3 of Stanley's book [66] here.

**Definition 2.6.1.** Let $P$ be a finite graded poset with a $\hat{0}$ and a $\hat{1}$. A map from the edges of the Hasse diagram of $P$ to the integers $\lambda : \mathcal{E}(P) \to \mathbb{Z}$ is an *R-labeling* for $P$ if for every interval $[x, y]$ of $P$, there is a unique saturated chain $C := x = x_1 \lessdot x_2 \lessdot \cdots \lessdot x_j = y$ where $\lambda(x, x_2) \leq \lambda(x_2, x_3) \leq \cdots \leq \lambda(x_{j-1}, y)$.

A notable difference between EL-labelings and R-labelings is the fact that we do not require the unique increasing chain in an interval to have a lexicographically smallest label sequence. Axel Hultman informed us in a personal communication that the proof of Theorem 2.4.4 may easily be modified to yield the following statement: Let $L$ be a finite lattice with set of join-irreducibles JoinIrred$(L)$ and $k = |\text{JoinIrred}(L)|$. Then the labeling $\lambda_\gamma$ induced by each choice of order-preserving bijection $\gamma : \text{JoinIrred}(L) \to [k]$ is an R-labeling if and only if $L$ is semimodular.

If $M = M(S)$ is a matroid (Definition 2.2.6) of rank $r$ on a finite set $S$, the *independence complex* of $M$ is the $(r-1)$-dimensional simplicial complex formed by the family of all independent sets in $M$. As we saw in the the discussion in Section 2.2, a geometric lattice is the lattice of flats, or closed sets, of a matroid.

*Remark.* There is a well-known result concerning matroid complexes which has a similar flavor to Theorem 2.4.4. This appears e.g. in the article "Homology and shellability of matroids and geometric lattices" of Björner [12] as Theorem 7.3.4. The statement of this result is as follows: a simplicial complex $\Delta$ is the independence complex of a matroid if and only if $\Delta$ is pure and every ordering of the vertices induces a shelling of $\Delta$. It seems interesting to note the resemblance between the necessary hypotheses for Theorem 7.3.4 of [12] and those of our characterization(s) of geometric lattices. It is natural to ask if one result is a translation of the other into a different language. This does not appear to be the case, rather the two results seem to be fundamentally quite distinct.

Finite geometric lattices, finite semimodular lattices, and finite supersolvable lattices are now known to have characterizations in terms of the types of edge-labelings they admit. At this point in time we do not know if there are other classes of lattices or posets that can be characterized in this fashion.

# Chapter 3

# Towards a Shellability Proof of a Hypergeometric Identity

## 3.1 Introduction

The goal of this chapter is to give results that could be employed towards a new proof of the identity

$$\sum_{s=0}^{n}(-1)^s\binom{n}{s}^3 = \begin{cases} 0 & \text{if } n \text{ is odd,} \\ (-1)^{n/2}\binom{3n/2}{n/2,n/2,n/2}, & \text{if } n \text{ is even,} \end{cases} \tag{3.1}$$

using tools from topological combinatorics. The identity (3.1) is due to Alfred Cardew Dixon (1865-1936). Dixon originally proved (3.1) in the paper "On the sum of the cubes of the coefficients in a certain expansion by the binomial theorem" in *Messenger of Mathematics* Volume 20 [30], which is a journal that ceased to publish in 1929. See page 121 of the book *Combinatory Analysis-Two Volumes in One* by Percy Alexander MacMahon (1854-1929) [47] for an application of (3.1).

This identity is actually a special case of a more general identity. Let $n_1, n_2$, and $n_3$ be nonnegative integers and let $N = n_1 + n_2 + n_3$. Then

$$\sum_{s=\max(0,n_1-n_2,n_3-n_2)}^{\min(n_1,n_3,n_1+n_3-n_2)}\binom{n_3}{s}\binom{n_2}{n_1-s}\binom{n_1}{n_2-n_3+s}(-1)^s = \\ \begin{cases} 0 \text{ if } N \text{ is odd,} \\ (-1)^{N/2-n_2}\binom{N/2}{N/2-n_1,N/2-n_2,N/2-n_3} \text{ if } N \text{ is even.} \end{cases} \tag{3.2}$$

The identity (3.1) is the case $n_1 = n_2 = n_3 = n$. The general case (3.2) is called the *well-poised $_3F_2$ transformation*. It is due to Wilfrid Norman Bailey (1893-1961) and can be found on page 97 of his book *Generalized Hypergeometric Series* [5]. See the proof of Lemma 4.2 in the paper [53] of Victor Reiner, Dennis Stanton, and Volkmar Welker for an application of the identity (3.2) related to the Charney-Davis conjecture. We thank Noam Elkies for suggesting the idea of proving the identity (3.1) using ideas from topology such as trying to interpret the alternating sum as an Euler characteristic. We thank Dennis Stanton for sharing his expertise and historical knowledge regarding the identities (3.1) and (3.2).

We now explain one approach to finding a topological proof of (3.1) that we work towards in this chapter. By Theorem 1.3.3, when a simplicial complex $\Delta$ is shellable (Definition 1.3.2), the Betti numbers $\beta_i(\Delta)$ (Definition 1.1.10) can be interpreted as counting the number of $i$-dimensional faces attaching to $\Delta$ along their entire boundary in a shelling order.

So, our aim is to find, for each $n$, a shellable simplicial complex $\Delta$ with face numbers (Definition 1.1.5) $f_{s-1} = \binom{n}{s}^3$ and then calculate the Betti numbers of $\Delta$. Then the Euler-Poincaré relation, which we saw as Equation 1.1 in Chapter 1:

$$\sum_{i=-1}^{d} (-1)^i f_i = \sum_{i=-1}^{d} (-1)^i \beta_i,$$

where $d$ is the maximum dimension of a face of $\Delta$, will hopefully give a new way of understanding and proving (3.1). Determining the Betti numbers may also lead to a refinement or a greater understanding of (3.1). Our candidate for a suitable simplicial complex to work towards these goals was given to us by Patricia Hersh, and is defined here:

**Definition 3.1.1** (Hersh). Fix $n \geq 1$. Let $\Delta(n)$ be the simplicial complex with vertices given by 3-tuples $(i_s, j_s, k_s)$ for $i_s, j_s, k_s \in [n]$ and faces given by collections of vertices

$$\{(i_1, j_1, k_1), \ldots, (i_r, j_r, k_r)\}$$

satisfying

$$i_1 < i_2 < \cdots < i_r \quad \text{and} \quad j_1 < j_2 < \cdots < j_r \quad \text{and} \quad k_1 < k_2 < \cdots < k_r.$$

It is easy to see that the number of $r$-faces of $\Delta(n)$ is counted by the product $\binom{n}{r+1}^3$. So

$$\overline{\chi}(\Delta(n)) = \sum_{s=0}^{n} (-1)^{s+1} \binom{n}{s}^3 = (-1) \times \sum_{s=0}^{n} (-1)^s \binom{n}{s}^3.$$

Therefore, if our approach to this problem works, the equation produced by writing down the Euler-Poincaré relation for $\Delta(n)$ will actually be equivalent to multiplying both sides of Equa-

tion 3.1 by $-1$. So, when we calculate the alternating sum of the Betti numbers $\sum_{i=-1}^{d}(-1)^i\beta_i$ (as we do for $n \leq 4$ in Section 3.4.1) we will actually check $(-1) \times \sum_{i=-1}^{d}(-1)^i\beta_i$ to verify (3.1).

We can generalize the construction of $\Delta(n)$. In particular, define a simplicial complex $\Gamma_p(n)$ with vertices given by sequences in $[n]^p$ for $p \geq 1$, and faces given by collections of vertices

$$\{(i_{1,1},\ldots,i_{1,p}),(i_{2,1},\ldots,i_{2,p}),\ldots,(i_{r,1},\ldots,i_{r,p})\}$$

satisfying $i_{\ell,a} \in [n]$ for all $\ell \in [r]$ and all $a \in [p]$, and $i_{\ell,a} < i_{(\ell+1),a}$ for all $\ell \in [r-1]$ and all $a \in [p]$. Then $\Gamma_p(n)$ has face numbers given by

$$f_{s-1} = \binom{n}{s}^p$$

for $0 \leq s \leq n$.

Note that for $p = 1$, $\Gamma_p(n)$ corresponds to the identity

$$\sum_{k=0}^{n}(-1)^k\binom{n}{k} = 0, \quad n \geq 1 \tag{3.3}$$

which appears as Exercise 1.3-(f) in *Enumerative Combinatorics Volume I* by Richard Stanley [66]. $\Gamma_1(n)$ is the traditional $n$-simplex $\Delta_{n-1}$ with vertices labeled with the labels $\{1,\ldots,n\}$ and therefore has the homotopy type (Definition 1.1.8) of a point. Therefore the left-hand side of (3.3) is $(-1)\widetilde{\chi}(\Delta_{n-1})$ which is equal to zero.

For $p = 2$, $\Gamma_p(n)$ corresponds to the identity

$$\sum_{k=0}^{n}(-1)^k\binom{n}{k}^2 = \begin{cases} 0 & \text{if } n \text{ is odd,} \\ (-1)^{n/2}\binom{n}{n/2}, & \text{if } n \text{ is even.} \end{cases} \tag{3.4}$$

Determining the value of the right-hand side of (3.4) appears as Exercise 5.48 in *A Course in Enumeration* by Martin Aigner [1]. One way to prove (3.4) is by the use of a sign-reversing involution (Definition 3.5.1).

Since we originally wanted to explore this new approach to the identity (3.1), we present results in this chapter for $p = 3$, i.e. for $\Delta(n)$. However, as the reader will see in Sections 3.3 and 3.4.1, the arguments leading to these results depend on the fact that $p$ is a positive integer and $p > 1$, but not truly on the fact that $p = 3$. So, the construction should generalize to create a family of non-pure, disconnected shellable simplicial complexes $\{\Gamma_p(n) : p \geq 1\}$. However, it is unknown if the Betti numbers of these generalized constructions for arbitrary $p$ will have any combinatorial significance of the type that the Betti numbers of $\Delta(n)$ do. We hope that the study begun in this chapter may eventually lead to the establishment of new identities.

We also note that we can view the simplicial complex $\Delta(n)$ as the order complex (Definition 1.1.3) of the poset $P_{\Delta(n)}$ where the elements of $P_{\Delta(n)}$ are integer triples in $[n]^3$ and $(i, j, k) \leq (i', j', k')$ in $P_{\Delta(n)}$ if and only if $i < i'$, $j < j'$ ,and $k < k'$ as integers. Our proof of the shellability of the complex $\Delta(n)$ does not directly use the fact that $\Delta(n)$ is an order complex of a poset, but we may incorporate this view into future work on this problem.

In Section 3.2 we observe some additional facts about $\Delta(n)$ and develop some technical language for operations on the faces of $\Delta(n)$. In Section 3.3 we establish a shelling order for the facets of $\Delta(n)$. In Section 3.4.1 we characterize the homology facets (Definition 3.4.1) of $\Delta(n)$, and use this characterization to verify our reformulation of (3.1) for $n \leq 4$. In Section 3.5 we discuss some possible continuations of the work described in this chapter.

## 3.2 Some Facts About $\Delta(n)$

When $n = 1$, the only nonempty face of $\Delta(n)$ is $\{(1, 1, 1)\}$. Figures 3.1 and 3.2 show the simplicial complexes $\Delta(2)$ and $\Delta(3)$, respectively. It is immediately apparent that $\Delta(n)$ is nonpure for all $n > 1$, and always has precisely one $(n-1)$-dimensional facet given by the collection of vertices $\{(1, 1, 1), (2, 2, 2), \ldots, (n, n, n)\}$. This is the maximum possible dimension of a face of $\Delta(n)$, so when we calculate the Euler-Poincaré relation it will be sufficient to calculate

$$\sum_{i=-1}^{n-1} (-1)^i f_i = \sum_{i=-1}^{n-1} (-1)^i \beta_i. \tag{3.5}$$

Note that each of $\Delta(1)$, $\Delta(2)$, and $\Delta(3)$ contain isolated vertices. Since $\Delta(1)$ is a point, it is pure and connected, but $\Delta(2)$ and $\Delta(3)$ are disconnected and nonpure.



Figure 3.1: The simplicial complex $\Delta(2)$.

Figure 3.2:  The simplicial complex $\Delta(3)$.

Since we will be establishing a shelling order for $\Delta(n)$, it is essential that we understand which faces are facets. The facets of $\Delta(n)$ are characterized in Lemma 3.2.1.

**Lemma 3.2.1.** *Let* $F = \{v_1, \ldots, v_r\} = \{(i_1, j_1, k_1), \ldots, (i_r, j_r, k_r)\}$ *be a face of* $\Delta(n)$. *Then* $F$ *is a facet if and only if* $F$ *satisfies the following three properties:*

*(P1)* $\max\{i_r, j_r, k_r\} = n$.

*(P2)* $\min\{i_1, j_1, k_1\} = 1$.

*(P3) If* $r \geq 2$, $\min\{i_{\ell+1} - i_\ell, j_{\ell+1} - j_\ell, k_{\ell+1} - k_\ell\} = 1$ *for all* $\ell \in [r-1]$.

*Proof.* Let $F = \{v_1, \ldots, v_r\}$ be a facet of $\Delta(n)$. Properties P1 and P2 clearly must hold for $F$: if P1 does not hold, then $F \subset \{v_0\} \cup F$ where $v_0 = (i_1 - 1, j_1 - 1, k_1 - 1)$. If P2 does not hold $F \subset F \cup \{v_{r+1}\}$, where $v_{r+1} = (i_r + 1, j_r + 1, k_r + 1)$. If P3 does not hold, there exists an index $\ell \in [r-1]$ where $\min\{i_{\ell+1} - i_\ell, j_{\ell+1} - j_\ell, k_{\ell+1} - k_\ell\} > 1$ and we can construct a vertex $v_s = (i_s, j_s, k_s)$ satisfying

$$i_\ell < i_s < i_{(\ell+1)}, \quad j_\ell < j_s < j_{(\ell+1)}, \quad \text{and} \quad k_\ell < k_s < k_{(\ell+1)}.$$

40

Then $F' = \{v_1, v_2, \ldots, v_\ell, v_s, v_{\ell+1}, \ldots, v_r\}$ properly contains $F$ and $F$ cannot be a facet. So each condition is necessary and sufficient for $F$ to be a facet. $\qquad \square$

From this characterization of the facets of $\Delta(n)$, we immediately obtain the next lemma:

**Lemma 3.2.2.** *For $n \geq 2$, the simplicial complex $\Delta(n)$ is nonpure and disconnected.*

*Proof.* Any vertex $v_s = (i_s, j_s, k_s)$ satisfying $\{1, n\} \subset \{i_s, j_s, k_s\}$ must be an isolated vertex, as $v_s$ cannot be contained in any other face in this case. If $n > 1$ there is more than one vertex in $\Delta(n)$, so $\Delta(n)$ contains isolated vertices and is disconnected for all $n \geq 2$. For all $n$ there is an $(n-1)$-dimensional facet $F(n) = \{(1,1,1), (2,2,2), \ldots, (n,n,n)\}$, and for $n \geq 2$, $\dim F(n) > 0$. So $\Delta(n)$ is not pure for $n \geq 2$. $\qquad \square$

It is also useful to obtain new faces of $\Delta(n)$ from old, and understand how to obtain a new facet from an old facet. To make these actions possible, we now define operations on the faces of $\Delta(n)$ in Definitions 3.2.3, 3.2.4, 3.2.6 and 3.2.9.

**Definition 3.2.3.** Let $F = \{v_1, v_2, \ldots, v_r\}$ be a face of $\Delta(n)$. For $\ell \in [r-1]$, let

$$A_\ell = \{a_\ell \in \{i_\ell, j_\ell, k_\ell\} : a_{\ell+1} - a_\ell > 1\}$$

and let

$$A_r = \{a_r \in \{i_r, j_r, k_r\} : a_r < n\}.$$

For $\ell$ such that $A_\ell \neq \emptyset$, define a **up-twist** $G$ about the vertex $v_\ell$ as the face of $\Delta(n)$ obtained by replacing $v_\ell$ in $F$ with the new vertex $v'_\ell$ of $\Delta(n)$ obtained by increasing an index of $v_\ell$ in $A_\ell$ by 1.

**Definition 3.2.4.** Let $F = \{v_1, v_2, \ldots, v_r\}$ be a face of $\Delta(n)$. For $\ell \in \{2, 3, \ldots, r\}$, let

$$B_\ell = \{b_\ell \in \{i_\ell, j_\ell, k_\ell\} : b_\ell - b_{\ell-1} > 1\}$$

and let

$$B_1 = \{b_1 \in \{i_1, j_1, k_1\} : b_1 > 1\}$$

For $\ell$ such that $B_\ell \neq \emptyset$, define a **down-twist** $G$ about the vertex $v_\ell$ as the face of $\Delta(n)$ obtained by replacing $v_\ell$ in $F$ with the new vertex $v'_\ell$ of $\Delta(n)$ obtained by decreasing an index of $v_\ell$ in $B_\ell$ by 1.

**Example 3.2.5.** In this example, we consider faces of $\Delta(5)$. In the face $F = \{(1,1,2), (2,5,5)\}$ with $v_1 = (1,1,2)$, and $v_2 = (2,5,5)$ the set $A_1 = \{j_1, k_1\} = \{1, 2\}$, and $B_2 = \{j_2, k_2\} = \{5, 5\}$. An up-twist of $F$ about $v_1$ is the new face $\{(1,2,2), (2,5,5)\}$. A down-twist of $F$ about $v_2$ is the new face $\{(1,2,2), (2,4,5)\}$.

**Definition 3.2.6.** Recall conditions P1, P2, and P3 from Lemma 3.2.1, and let $F$ be a facet. We say an up-twist $G$ of a facet $F = \{v_1, v_2, \ldots, v_r\}$ is **safe** if either $\ell > 1$ and condition P3 is conserved, or $\ell = 1$ and condition P2 is conserved. We say a down-twist $G$ of a facet $F$ is safe if either $\ell < r$ and condition P3 is conserved, or $\ell = r$ and condition P1 is conserved.

**Example 3.2.7.** In this example, we again consider faces of $\Delta(5)$. The down-twist

$$\{(1, 2, 2), (2, 4, 4)\}$$

of

$$\{(1, 2, 2), (2, 4, 5)\}$$

about the vertex $(2, 4, 5)$ is not safe. The down-twist

$$\{(1, 2, 2), (2, 3, 5)\}$$

of

$$\{(1, 2, 2), (2, 4, 5)\}$$

about the vertex $(2, 4, 5)$ is safe.

**Lemma 3.2.8.** *Let $F_1 = \{v_1, \ldots, v_r\}$ be a facet of $\Delta(n)$. If $F_2$ is a safe up-twist or a safe down-twist of $F_1$, then then $F_2$ is a facet of $\Delta(n)$.*

*Proof.* First we consider the case where $F_1$ is a facet $\{v_1, \ldots, v_r\}$ of $\Delta(n)$ and

$$F_2 = \{v_{\ell'}\} \cup \{v_j : j \in [r] \setminus \{\ell\}\}$$

is a safe up-twist of $F_1$ about the vertex $v_\ell$. By Lemma 3.2.1, it is sufficient to show that $F_2$ satisfies P1, P2, and P3. If $\ell < r$, then $v_r$ is unaffected by the vertex change and P1 holds for $F_2$. If $\ell = r$, then the maximum element of $v'_r$ will not decrease and P1 is satisfied by $F_2$. If $\ell > 1$, then P2 is trivially satisfied by $F_2$. If $\ell = 1$, then since $F_2$ is a *safe* up-twist, $\min\{i_1, j_1, k_1\} = \min\{i'_1, j'_1, k'_1\} = 1$, where $(i'_1, j'_1, k'_1)$ is the new vertex in $v'_1 \in F_2$. So $P2$ holds.

Now, since $F_1$ is a facet, $\min\{i_{\ell+1} - i_\ell, j_{\ell+1} - j_\ell, k_{\ell+1} - k_\ell\} = 1$. Without loss of generality, we can say $\min\{i_{\ell+1} - i_\ell, j_{\ell+1} - j_\ell, k_{\ell+1} - k_\ell\} = i_{\ell+1} - i_\ell$, so that $i_\ell \notin A_\ell$ and $v'_\ell = (i_\ell, j'_\ell, k'_\ell)$. Therefore $\min\{i_{\ell+1} - i_\ell, j_{\ell+1} - j'_\ell, k_{\ell+1} - k'_\ell\} = 1$ and P3 holds.

A similar argument shows that if $F_2$ is a safe down-twist of $F_1$, then $F_2$ is also a facet of $\Delta(n)$. $\square$

**Definition 3.2.9.** Let $F = \{v_1, \ldots, v_r\}$ be a face of $\Delta(n)$ where $r \geq 2$. For $\ell \in [r-1]$, a **contraction** of $F$ about $v_\ell$ is a face

$$G = \{w\} \cup \{v_j : j \in [r] \setminus \{\ell, \ell+1\}\}$$

where $w = (i, j, k)$ is a vertex satisfying

(ii) $i \in \{i_\ell, i_{\ell+1}\}, j \in \{j_\ell, j_{\ell+1}\}, k \in \{k_\ell, k_{\ell+1}\}$,

(ii) $\min\{i_{\ell+2} - i, j_{\ell+2} - j, k_{\ell+2} - k\} = 1$, and

(iii) $\min\{i - i_{\ell-1}, j - j_{\ell-1}, k - k_{\ell-1}\} = 1$.

We say a contraction is **safe** if $F$ is a facet and either $\ell = 1$ and $\min\{i, j, k\} = 1$ or $\ell + 1 = r$ and $\max\{i, j, k\} = n$.

**Example 3.2.10.** In this example, we consider faces in $\Delta(6)$. If

$$F = \{(1, 2, 1), (2, 4, 3), (3, 5, 5), (6, 6, 6)\},$$

then $G_1 = \{(1, 2, 1), (2, 5, 5), (6, 6, 6)\}$ is obtained from a contraction of $F$ about $v_2 = (2, 4, 3)$, and $G_2 = \{(1, 4, 3), (3, 5, 5), (6, 6, 6)\}$ and $G_3 = \{(2, 2, 3), (3, 5, 5), (6, 6, 6)\}$ are obtained from contractions of $F$ about $v_1 = (1, 2, 1)$. Note that $G_2$ is obtained from a safe contraction, but $G_3$ is not.

The next lemma shows that we can use safe contractions to obtain a new facet from an old facet in $\Delta(n)$.

**Lemma 3.2.11.** *Let $F_1 = \{v_1, \ldots, v_r\}$ be a facet of $\Delta(n)$ and let $\ell \in [r-1]$. If*

$$F_2 = \{w\} \cup \{v_j : j \in [r] \setminus \{\ell, \ell+1\}\}$$

*where $w = (i, j, k)$ is obtained from $F_1$ via a safe contraction about $v_\ell$, then $F_2$ is a facet of $\Delta(n)$.*

*Proof.* We appeal again to Lemma 3.2.1; we must show that $F_2$ satisfies P1, P2, and P3. If $\ell < r - 1$, then $v_r \in F_2$ and $F_2$ inherits P1 from $F_2$. If $\ell = r - 1$, then $\ell + 1 = r$ and since $F_2$ is obtained from a safe contraction, $\max\{i, j, k\} = \max\{i_r, j_r, k_r\} = n$. So $F_2$ satisfies P1. If $\ell > 1$, then $v_1 \in F_2$ and $F_2$ automatically inherits P2 from $F_1$. If $\ell = 1$, then since $F_2$ is obtained from a safe contraction, $\min\{i, j, k\} = \min\{i_1, j_1, k_1\} = 1$ and $F_2$ satisfies P2.

Conditions $(ii)$ and $(iii)$ in the definition of a contraction guarantee that $F_2$ satisfy P3. So P1, P2, and P3 are true for $F_2$, and $F_2$ is a facet of $\Delta(n)$. $\qquad\qquad\square$

## 3.3   A Shelling Order for $\Delta(n)$

In this section we construct a shelling order for $\Delta(n)$. Recall that by Lemma 3.2.2, $\Delta(n)$ is not pure. To set up the shelling order, we first partition $\Delta(n)$ into sets of facets according to dimension. For $0 \leq m \leq n-1$, let $S_m$ be the set of facets of dimension $m$. For example, $S_{n-1}$ is the set containing the single $(n-1)$-dimensional facet $F(n) = \{(1,1,1),(2,2,2),(3,3,3),...,(n,n,n)\}$, and $S_0$ contains the facets comprised solely of isolated vertices such as the facet $F = \{(1,n,1)\}$. The next definition allows us to order the facets in $S_m$ for a fixed $m$ using the lexicographic order.

**Definition 3.3.1.** The $\sigma$-*word* $\sigma(F)$ of the face $F = \{v_1, \ldots, v_r\}$ is the sequence

$$(i_1, j_1, k_1, i_2, j_2, k_2, i_3, \ldots, k_{r-1}, i_r, j_r, k_r).$$

Informally, we can see that the $\sigma$-word of a face is obtained by simply ignoring all the parentheses in the listing of the vertices of the face.

**Example 3.3.2.** The $\sigma$-word of the vertices of the facet

$$F = \{(1,2,1),(3,3,3),(4,4,5)\}$$

of $\Delta(5)$ is
$$\sigma(F) = (1,2,1,3,3,3,4,4,5).$$

The $\sigma$-word of the vertices of the facet

$$G = \{(1,2,2),(2,3,3),(4,4,5)\}$$

of $\Delta(5)$ is
$$\sigma(G) = (1,2,2,2,3,3,4,4,5).$$

In the lexicographic order we have $\sigma(F) < \sigma(G)$.

Now we define the order on the facets of $\Delta(n)$ that we will show is a shelling order.

**Definition 3.3.3.** Define an order $\mathcal{O}$ on the facets of $\Delta(n)$ as folllows: $F_i < F_j$ in the order $\mathcal{O}$ if

1. $F_i \in S_m$ and $F_j \in S_\ell$ for $\ell < m$ or

2. $\ell = m$ and $\sigma(F_i)$ is lexicographically smaller than $\sigma(F_j)$.

**Theorem 3.3.4.** *The order $\mathcal{O}$ is a shelling order for the facets of $\Delta(n)$.*

We need the following well-known lemma, which provides a useful working definition of a shelling, in our proof of Theorem 3.3.4. This lemma is explicitly stated for non-pure simplicial complexes as Lemma 2.3 in the paper [15] of Anders Björner and Michelle Wachs, but a version for pure simplicial complexes appears earlier in the paper [11] of Björner.

**Lemma 3.3.5.** *An order $F_1, F_2, \ldots, F_t$ of the facets of a simplicial complex $\Delta$ is a shelling if and only if for every $i$ and $k$ satisfying $1 \leq i < k \leq t$ there is a $j$ with $1 \leq j < k$ and a vertex $v \in F_k$ such that $F_i \cap F_k \subset F_j \cap F_k = F_k \setminus \{v\}$.*

We will follow the notation of Lemma 3.3.5 and let $[t]$ denote the index set for the order $\mathcal{O}$. To work with Lemma 3.3.5 in the proof of Theorem 3.3.4, we will be fixing two facets $F_i$ and $F_k$ and constructing a facet $F_j$ satisfying the conditions of the lemma. To make this easier, we now develop some notation for vertex subsets of $F_i$ and $F_k$.

Let $F_k$, for $k > 1$, be a facet of $\Delta(n)$. Let $i$ be an index satisfying $1 \leq i < k \leq t$. Let $\mathcal{V}_{i,k}$ denote the (possibly empty) set of vertices in $F_i \cap F_k$, let $\mathcal{V}_k = F_k \setminus F_i$, and let $\mathcal{V}_i = F_i \setminus F_k$. Write $\mathcal{V}_{i,k} = \{v_{c,1}, \ldots, v_{c,s}\}$, $\mathcal{V}_k = \{v_{k,1}, \ldots, v_{k,e}\}$, and $\mathcal{V}_i = \{v_{i,1}, \ldots, v_{i,u}\}$. We write the vertex sets so that as positive integers, $(c, 1) < \cdots < (c, s)$, $(k, 1) < \cdots < (k, e)$, and $(i, 1) < \cdots < (i, u)$. Also, we write the indices of $\mathcal{V}_{i,k}$, $\mathcal{V}_k$, and $\mathcal{V}_i$ in the same order they appear in $F_k$ and $F_i$, and we do not rename the indices when considering the subsets $\mathcal{V}_{i,k}, \mathcal{V}_i$, and $\mathcal{V}_k$.

**Example 3.3.6.** Let
$$F_i = \{(1, 2, 1), (3, 3, 3), (4, 4, 4), (5, 5, 5)\}$$
and
$$F_k = \{(1, 2, 1), (2, 3, 3), (5, 4, 5)\}.$$

Then $\mathcal{V}_{i,k} = \{(1, 2, 1)\}$, $\mathcal{V}_k = \{(2, 3, 3), (5, 4, 5)\}$, and $\mathcal{V}_i = \{(3, 3, 3), (4, 4, 4), (5, 5, 5)\}$. Also, $\{(c, 1)\} = \{1\}$, $\{(k, 1), (k, 2)\} = \{2, 3\}$ and $\{(i, 1), (i, 2), (i, 3)\} = \{2, 3, 4\}$.

The next lemma will make it easier to work with the sets $\mathcal{V}_{i,k}$, $\mathcal{V}_i$, and $\mathcal{V}_k$.

**Lemma 3.3.7.** *Let $F_i$ and $F_k$ be facets of $\triangle(n)$ such that $i < k \in [t]$. There exist ordered partitions of $\mathcal{V}_{i,k}, \mathcal{V}_i$, and $\mathcal{V}_k$ into blocks of ordered vertices*

$$\mathcal{V}_{i,k} = C_1 | \cdots | C_N,$$

$$\mathcal{V}_i = I_1 | \cdots | I_M,$$

*and*

$$\mathcal{V}_k = K_1 | \cdots | K_M$$

*such that each ordered block of the ordered partitions corresponds to a consecutive subsequence of vertices in a facet.*

*Proof.* We can generate the required partitions of the vertices of $\mathcal{V}_i$, $\mathcal{V}_k$, and $\mathcal{V}_{i,k}$ using an algorithmic approach. We will explain the algorithm $\mathcal{V}_{i,k} = C_1|\cdots|C_N$; the algorithms for $\mathcal{V}_i$ and $\mathcal{V}_k$ are similar.

Let $F_i = \{v_1, \ldots, v_r\}$. If $\mathcal{V}_{i,k} = \emptyset$, then the partition is empty, and there is nothing to compute. So, assume $\mathcal{V}_{i,k} \neq \emptyset$. We use the following algorithm to build the blocks of the ordered partition $C_1|\cdots|C_N$.

**Algorithm 3.3.8.**    • Input: The vertices $\{v_1, \ldots, v_r\}$ of $F_i$ and the vertex subset $\mathcal{V}_{i,k}$.

- Output: An ordered partition of $\mathcal{V}_{i,k}$ of ordered blocks of vertices in $\mathcal{V}_{i,k}$, in which each block is a set of vertices that are both consecutive in $\{v_1, \ldots, v_r\}$ and written in the order that they appear in $\{v_1, \ldots, v_r\}$.

- Initialize $\ell(1) = \min\{\ell \in [r] : v_\ell \in \mathcal{V}_{i,k}\}$, set $C_1 = \{v_{\ell(1)}\}$ .

- While $\ell(i) < r$:

    - If $v_{\ell(i)+1} \in \mathcal{V}_{i,k}$, set $C_i = C_i \cup \{v_{\ell(i)+1}\}$, update $\ell(i) = \ell(i) + 1$.
    - Else if $v_{\ell(i)+1} \notin \mathcal{V}_{i,k}$, set $v_{\ell(i)}$ as the last vertex in $C_i$.
        * If the set $\{m \in [r] : m > \ell(i) \text{ and } v_m \in \mathcal{V}_{i,k}\}$ is empty, $C_i = C_N$ and the algorithm terminates.
        * Else update $\ell(i+1) = \min\{m \in [r] : m > \ell(i) \text{ and } v_m \in \mathcal{V}_{i,k}\}$ and set $C_{i+1} = \{v_{\ell(i+1)}\}$.

- Return the partition $\mathcal{V}_{i,k} = C_1|\cdots|C_N$.

It remains to explain why the partitions of $\mathcal{V}_i$ and $\mathcal{V}_k$ both have the same number of blocks $M$. Assume by way of contradiction that the partition of $F_i$ has more blocks than the partition of $F_k$. Then either $(i)$ there is at least one vertex in the sequence $\{v_1, \ldots, v_r\}$ that is between two vertices of $F_k$, $(ii)$ there is a vertex of $F_i$ greater than the last vertex of $F_k$, or $(iii)$ there is a vertex smaller than the first vertex of $F_k$. Case $(i)$ implies $F_k$ does not satisfy P3, Case $(ii)$ implies that $F_k$ does not satisfy P1, and Case $(iii)$ implies $F_k$ does not satisfy P2. So, all three cases are impossible by Lemma 3.2.1. So $F_i$ cannot have more blocks in the ordered partition than $F_k$. The argument is symmetric in $F_i$ and $F_k$, so the ordered partitions of $F_i$ and $F_k$ have the same number of blocks. □

**Example 3.3.9.** Here is an example of the ordered partitions with ordered blocks described in Lemma 3.3.7. Let

$$F_i = \{(1,2,1), (2,3,3), (3,4,4), (5,5,5), (6,6,6), (7,8,8), (9,9,9), (10,10,10)\}$$

and

$$F_k = \{(1,2,1), (2,3,5), (6,6,6), (7,8,9), (10,10,10)\}$$

be facets of $\Delta(10)$. Then

$$\mathcal{V}_{i,k} = \{(1,2,1), (6,6,6), (10,10,10)\},$$

$$\mathcal{V}_i = \{(2,3,3), (3,4,4), (5,5,5), (7,8,8), (9,9,9)\},$$

and

$$\mathcal{V}_k = \{(2,3,5), (7,8,9)\}.$$

We have

$$\mathcal{V}_{i,k} = C_1|C_2|C_3 = \{(1,2,1)\}|\{(6,6,6)|\{(10,10,10)\},$$

$$\mathcal{V}_i = I_1|I_2 = \{(2,3,3), (3,4,4), (5,5,5),\}|\{(7,8,8), (9,9,9)\},$$

and

$$\mathcal{V}_k = K_1|K_2 = \{(2,3,5)\}|\{(7,8,9)\}.$$

The next lemma is useful in our proof of Theorem 3.3.4. Recall that we write $\mathcal{V}_k = \{v_{k,1}, \ldots, v_{k,e}\}$.

**Lemma 3.3.10.** *Let $i < k$ be indices in the order $\mathcal{O}$. There exists $\ell \in \{(k,1)\ldots,(k,e)\}$ such that $B_\ell \neq \emptyset$.*

*Proof.* First we handle the case where $\dim F_i = \dim F_k$. In this case $F_i$ and $F_k$ each have $r$ vertices, and the sequences $\sigma(F_i)$ and $\sigma(F_k)$ are both of length $3r$. By our construction of the shelling order $\mathcal{O}$ in Definition 3.3.3 this implies $\sigma(F_i) < \sigma(F_k)$ in the lexicographic order which means the first place the two sequences differ, call this index $b \in [3r]$, is larger in $\sigma(F_k)$.

In other words, we have

$$\sigma(F_i) = (p_1, \ldots, p_{3r})$$

and

$$\sigma(F_k) = (q_1, \ldots, q_{3r})$$

where $p_a = q_a$ for $a \in [3r]$ satisfying $a < b$, and $q_b > p_b$ as integers.

The first place the sequences differ occurs in the vertex of smallest index not present in both

47

$F_i$ and $F_k$. So $q_b \in \{i_{k,1}, j_{k,1}, k_{k,1}\}$ as $v_{k,1}$ is the vertex of smallest index in $\mathcal{V}_k$. Without loss of generality, we can say $b$ designates the position of $i_{k,1}$. Recall that we write $\mathcal{V}_i = \{v_{i,1}, \ldots, v_{i,u}\}$. Then we have $i_{k,1} > i_{i,1}$. If $(k,1) = (i,1) = 1$, then $i_{k,1} \geq 2$ and $B_{k,1} \neq \emptyset$. If $(k,1) > 1$, $i_{(k,1)-1}$ must appear in a vertex in $\mathcal{V}_{i,k}$ by our choice of $b$, and we can write $i_{k,1} - i_{(k,1)-1} > i_{i,1} - i_{(k,1)-1} \geq 1$. So $B_{k,1} \neq \emptyset$ in this case.

Next we handle the case where $\dim F_i > \dim F_k$. We can write $F_i$ and $F_k$ as the disjoint unions

$$F_k = \mathcal{V}_{i,k} \sqcup \mathcal{V}_k, \quad F_i = \mathcal{V}_{i,k} \sqcup \mathcal{V}_i.$$

We know that $|\mathcal{V}_i| > |\mathcal{V}_k|$ because $\dim F_i > \dim F_k$. By Lemma 3.3.7 there exist partitions of $\mathcal{V}_{i,k}, \mathcal{V}_i$, and $\mathcal{V}_k$

$$\mathcal{V}_{i,k} = C_1 | \cdots | C_N, \quad \mathcal{V}_i = I_1 | \cdots | I_M,$$

and

$$\mathcal{V}_k = K_1 | \cdots | K_M$$

such that each block in each partition corresponds to an uninterrupted sequence of vertices in a facet. Since $|\mathcal{V}_i| > |\mathcal{V}_k|$ and the ordered partitions of $\mathcal{V}_i$ and $\mathcal{V}_k$ have the same number of blocks, there must exist $A \in [M]$ such that $|I_A| > |K_A|$. For indices $x \in [u]$ and $y \in [e]$ we can write

$$I_A = \{v_{i,x}, \ldots, v_{i,(x+|I_A|)}\}, \quad \text{and} \quad K_A = \{v_{k,y}, \ldots, v_{k,(y+|K_A|)}\}.$$

Recall that the $\sigma$-word $\sigma(K_A)$ (Definition 3.3.1) is the ordered set of indices of all vertices appearing in the face $K_A$. We divide the proof for $\dim F_i > \dim F_k$ into two cases

$$n \in \sigma(K_A), \tag{3.6}$$

and

$$n \notin \sigma(K_A). \tag{3.7}$$

Consider first the case (3.6). This case implies $n \in \{i_{k,(y+|K_M|)}, j_{k,(y+|K_M|)}, k_{k,(y+|K_M|)}\}$, the index set of the last vertex in the last block of the partition $K_1 | \cdots | K_M$ of $\mathcal{V}_k$. We can assume without loss of generality that $n = i_{k,(y+|K_M|)}$. Then $n$ appears as an element of $v_{i,(x+|I_M|)}$ also. Because of this, we know that in $F_i$, the vertices in $I_M$ immediately follow the vertices in $C_N$, and in $F_k$, the vertices in $K_M$ immediately follow the vertices in $C_N$. For some $z \in [s]$, we can write $C_N = \{v_{c,z}, \ldots, v_{c,(z+|C_N|)}\}$.

Then $n - i_{c,(z+|C_N|)} \geq |I_M|$, and the net change in the $i$ index in the vertices of $K_M$ is bounded below by $|I_M| > |K_M|$, and there are only $|K_M|$ vertices to accomplish this change. Therefore there must exist an index $\ell \in \{(k,y), \ldots, (k, y + |K_M|)\}$ such that $i_\ell - i_{\ell-1} > 1$. So,

for this $v_\ell \in \mathcal{V}_k$, $B_\ell \neq \emptyset$.

Now we consider the case (3.7). This implies that for some $w \in [s]$ there exists a vertex $v_{c,w} \in \mathcal{V}_{i,k}$ where $(c,w) = (k, (y+|K_A|))+1$ in the vertex numbering in $F_k$. Since $F_k$ is a facet, it satisfies P3 from Lemma 3.2.1, which means that $\min\{i_{c,w} - i_{k,(y+|K_A|)}, j_{c,w} - j_{k,(y+|K_A|)}, k_{c,w} - k_{k,(y+|K_A|)}\} = 1$. Without loss of generality we can say $i_{c,w} - i_{k,(y+|K_A|)} = 1$. If $1 \in \sigma(K_A)$ then $A = 1$ and $i_{k,(y+|K_1|)} \geq |I_1|$ where $|I_1| > |K_1|$, but we only have $|K_1|$ vertices to accomplish this index change and so there exists $\ell \in \{1, \ldots, (1+|K_1|)\}$ such that $B_\ell \neq \emptyset$.

If $1 \notin \sigma(K_A)$, there exists $x \in [s]$ and $v_{c,x} \in \mathcal{V}_{i,k}$ such that $(c,x)+1 = (k,y)$ in the label sequence of the vertices of $F_k$. Then $i_{k,(y+|K_A|)} - i_{c,x} \geq |I_A|$ where $|I_A| > |K_A|$. But we only have $|K_A|$ vertices to accomplish this index change and so there exists $\ell \in \{(k,y), \ldots, (k, (y+|K_A|))\}$ such that $B_\ell \neq \emptyset$. This completes the proof of the Lemma for the case $\dim F_i > \dim F_k$. $\qquad\square$

Now we proceed to the proof of Theorem 3.3.4. The essence of the proof is that given any pair of facets $F_i$ and $F_k$ such that $i < k$ in $\mathcal{O}$, we may use Lemma 3.3.10 to construct a facet $F_j$ such that the hypotheses of Lemma 3.3.5 is satisfied, which will show that $\mathcal{O}$ is a shelling order.

*Proof.* Let $F_i$ and $F_k$ be such that $i$ and $k$ satisfy $1 \leq i < k \leq t$ in the order $\mathcal{O}$. Recall we write $\mathcal{V}_{i,k} = F_i \cap F_k$ and $\mathcal{V}_k = \{v_{k,1}, \ldots, v_{k,e}\}$, where $\mathcal{V}_k = F_k \setminus F_i$. Write $F_k = \{v_1, \ldots, v_r\}$. We will find a vertex $v \in F_k$ and construct a facet $F_j$ such that $1 \leq j < k$ and such that $\mathcal{V}_{i,k} \subset F_j \cap F_k = F_k \setminus \{v\}$. This will show that $\mathcal{O}$ is a shelling order by Lemma 3.3.5.

By Lemma 3.3.10 there exists $\ell \in \{(k,1), \ldots, (k,e)\}$ such that $B_\ell \neq \emptyset$. We will divide the proof into two cases: $\ell = r$ and $\ell < r$. For now assume that $\ell < r$. If such an $\ell$ exists we choose $\ell$ that is minimal.

Then choose the "left-most" vertex element in $B_\ell$: for example if $B_\ell = \{i_\ell, k_\ell\}$ we choose $i_\ell$. Without loss of generality we can say that $i_\ell$ is the left-most element of the set $B_\ell$. Let $w = (i_\ell - 1, j_\ell, k_\ell)$. Since $F_k$ is a facet, we know that $\min\{i_{\ell+1} - i_\ell, j_{\ell+1} - j_\ell, k_{\ell+1} - k_\ell\} = 1$. We now have two sub-cases to consider: (i): $\min\{j_{\ell+1} - j_\ell, k_{\ell+1} - k_\ell\} = 1$ and (ii): $\min\{j_{\ell+1} - j_\ell, k_{\ell+1} - k_\ell\} > 1$. In the case (i), the down-twist (Definition 3.2.4) about $v_\ell$

$$F_j = F_k \setminus \{v_\ell\} \cup \{w\}$$

is safe and $F_j$ is a facet.

In this instance $\dim F_j = \dim F_k$. The only place $\sigma(F_j)$ and $\sigma(F_k)$ differ is the position of $i_\ell - 1$ from the new vertex $w$. So $\sigma(F_j) < \sigma(F_k)$ and we know $j < k$ in the order $\mathcal{O}$. Also, $\mathcal{V}_{i,k} \subset F_j \cap F_k = F_k \setminus \{v_\ell\}$, so $v_\ell$ and $F_j$ satisfy the conditions of Lemma 3.3.5.

Next, consider the sub-case (ii): $\min\{j_{\ell+1} - j_\ell, k_{\ell+1} - k_\ell\} > 1$. Since $F_k$ is a facet, P3 is

satisfied and $i_{\ell+1} - i_\ell = 1$ must hold. In this case the face

$$F_j = F_k \setminus \{v_\ell\} \cup \{w, (i_\ell, j_\ell + 1, k_\ell + 1)\}$$

satisfies P3 and is a facet. Since $\dim F_j > \dim F_k$, $j < k$ in $\mathcal{O}$. Clearly $\mathcal{V}_{i,k} \subset F_j \cap F_k = F_k \setminus \{v_\ell\}$.

Next we consider the case where the only $\ell \in \{(k, 1), \dots, (k, e)\}$ satisfying $B_\ell \neq \emptyset$ is $\ell = r$. Again without loss of generality we can say that $i_r$ is the left-most index in $B_r$. Let $w = (i_r - 1, j_r, k_r)$. There are two sub-cases to consider: (i) $n \in \{j_r, k_r\}$ and (ii) $n \notin \{j_r, k_r\}$.

If (i) $n \in \{j_r, k_r\}$, then the down-twist about $v_r$

$$F_j = F_k \setminus \{v_r\} \cup \{w\}$$

is safe and $F_j$ is a facet of the same dimension as $F_k$ satisfying $\sigma(F_j) < \sigma(F_k)$ and so $j < k$ in $\mathcal{O}$. For the sub-case (ii) when $n \notin \{j_r, k_r\}$, let

$$F_j = F_k \setminus \{v_r\} \cup \{w, (n, n, n)\}.$$

Since $\dim F_j > \dim F_k$, we have $j < k$ in $\mathcal{O}$. In both sub-cases $\mathcal{V}_{i,k} \subset F_j \cap F_k = F_k \setminus \{v_r\}$. This completes the proof. $\square$

## 3.4 The Homology Facets of $\Delta(n)$

One approach to calculating the Betti numbers (Definition 1.1.10) of $\Delta(n)$ is by understanding how the shelling order $\mathcal{O}$ puts $\Delta(n)$ together as a topological space. This helps us understand the topology of $\Delta(n)$. Towards this goal, we establish the next definition:

**Definition 3.4.1.** We say a facet $F_k$ of a simplicial complex $\Delta$ with shelling order $\mathcal{O}$ is a *homology $(r-1)$-facet* if $F_k$ is an $(r-1)$-dimensional facet satisfying

$$\partial F_k = F_k \cap \bigcup_{i<k} \overline{F_i}$$

where $i < k$ in $\mathcal{O}$, and $\partial F_k$ denotes the boundary complex of $F_k$, which is the sub-complex of $\Delta$ formed by taking the collection of all proper faces of $F_k$.

In other words, $F_k$ is a homology $(r-1)$-facet when $F_k$ attaches to $\Delta$ along its whole boundary in the shelling order $\mathcal{O}$. As we saw in Chapter 1, the Betti numbers of any shellable simplicial complex have a natural interpretation in terms of homology facets. By Theorem 1.3.3, the number of $(r-1)$-spheres in the homotopy type of $\Delta$ is the number of homology

$(r-1)$-facets. The next lemma characterizes the homology facets of $\Delta(n)$ for dimension 1 and greater.

**Lemma 3.4.2.** *Let $r \geq 2$. A facet $F_k = \{v_1, \ldots, v_r\}$ is a homology $(r-1)$-facet of $\Delta(n)$ if and only if $B_\ell \neq \emptyset$ for all $\ell \in [r]$.*

*Proof.* First let $B_\ell \neq \emptyset$ for all $\ell \in [r]$. It suffices to show that for all $\ell$, $F_k \setminus \{v_\ell\} \subset F_{j(\ell)}$ for some $j(\ell) < k$. First, let $\ell = r$. If at least two of the elements of the set $\{i_r, j_r, k_r\}$ are equal to $n$, then since $B_r \neq \emptyset$, we can say without loss of generality that $i_r - i_{r-1} > 1$. Then the facet

$$F_{j(\ell)} = F_k \setminus \{v_r\} \cup \{(i_r - 1, j_r, k_r)\}$$

satisfies $\dim F_{j(\ell)} = \dim F_k$ and $\sigma(F_{j(\ell)}) < \sigma(F_k)$, so $j(\ell) < k$ and we have the desired containment $F_k \setminus \{v_\ell\} \subset F_{j(\ell)}$. If $B_r = \{n\}$, without loss of generality we can say that $B_r = \{i_r\}$. Then $\min\{j_r, k_r\} = n - p$ for some $p \geq 1$. Let

$$F_{j(\ell)} = F_k \setminus \{v_r\} \cup \{(n-p, n-p, n-p), (n-p+1, n-p+1, n-p+1), \ldots, (n,n,n)\}.$$

Then $\dim F_{j(\ell)} > \dim F_k$, so $j(\ell) < k$ and $F_k \setminus \{v_r\} \subset F_{j(\ell)}$.

Now, let $\ell < r$. Either $|B_\ell| = 1$ or $|B_\ell| = 2$. (Since $F_k$ is a facet and satisfies P3, $|B_\ell| < 3$). If $|B_\ell| = 2$ we can assume $B_\ell = \{i_\ell, j_\ell\}$. If $\min\{j_{\ell+1} - j_\ell, k_{\ell+1} - k_\ell\} = 1$, then let

$$F_{j(\ell)} = F_k \setminus \{v_\ell\} \cup \{(i_\ell - 1, j_\ell, k_\ell)\}.$$

Then $\sigma(F_{j(\ell)}) < \sigma(F_k)$, with $F_k \setminus \{v_\ell\} \subset F_{j(\ell)}$ and $\dim F_k = \dim F_{j(\ell)}$, so $j(\ell) < k$. If $\min\{j_{\ell+1} - j_\ell, k_{\ell+1} - k_\ell\} > 1$ then $i_{\ell+1} - i_\ell = 1$ because $F_k$ is a facet and satisfies P3. Then let

$$F_{j(\ell)} = F_k \setminus \{v_\ell\} \cup \{(i_\ell - 1, j_\ell, k_\ell), (i_\ell, j_\ell + 1, k_\ell + 1)\}$$

and again $\dim F_{j(\ell)} > \dim F_k$, so $j(\ell) < k$ and $F_k \setminus \{v_\ell\} \subset F_{j(\ell)}$.

If $|B_\ell| = 1$, then we can assume $B_\ell = \{i_\ell\}$. If $i_{\ell+1} - i_\ell > 1$, then

$$F_{j(\ell)} = F_k \setminus \{v_\ell\} \cup \{(i_\ell - 1, j_\ell, k_\ell)\}$$

satisfies P3 (because $F_k$ does), $\dim F_k = \dim F_{j(\ell)}$ and $\sigma(F_{j(\ell)}) < \sigma(F_k)$, so $j(\ell) < k$ and $F_k \setminus \{v_\ell\} \subset F_{j(\ell)}$. If $i_{\ell+1} - i_\ell = 1$ and $\min\{j_{\ell+1} - j_\ell, k_{\ell+1} - k_\ell\} > 1$, then

$$F_{j(\ell)} = F_k \setminus \{v_\ell\} \cup \{(i_\ell - 1, j_\ell - 1, k_\ell - 1), (i_\ell, j_\ell, k_\ell)\}$$

satisfies $\dim F_k < \dim F_{j(\ell)}$ so $j(\ell) < k$ and $F_k \setminus \{v_\ell\} \subset F_{j(\ell)}$. So whenever $B_\ell \neq \emptyset$ for all $\ell \in [r]$,

$F_k$ attaches along its entire boundary in the shelling order $\mathcal{O}$ and is a homology facet.

For the converse, assume that $F_k = \{v_1, \ldots, v_r\}$ is a homology facet. We wish to show that $B_\ell \neq \emptyset$ for all $\ell \in [r]$. Assume by way of contradiction that there exists $\ell \in [r]$ where $B_\ell = \emptyset$. Since $F_k$ is a homology facet, $F_k \setminus \{v_\ell\} \subset F_{j(\ell)}$ for some $j(\ell) < k$. If $\dim F_k = \dim F_{j(\ell)}$, then

$$F_{j(\ell)} = F_k \setminus \{v_\ell\} \cup \{v'_\ell\}$$

for some $v'_\ell \neq v_\ell$, and $\sigma(F_{j(\ell)}) < \sigma(F_k)$. Then since the only entries in the sequences $\sigma(F_{j(\ell)})$ and $\sigma(F_k)$ that are different come from $v_\ell$ and $v'_\ell$, one of the three inequalities (i) $i'_\ell < i_\ell$, (ii) $j'_\ell < j_\ell$, or (iii) $k'_\ell < k_\ell$ must be true. If $i'_\ell < i_\ell = i_{\ell-1} + 1$, then this is a contradiction because $i'_\ell > i_{\ell-1}$. The same contradiction arises if inequalities (ii) or (iii) hold.

If $\dim F_{j(\ell)} > \dim F_k$, then

$$F_{j(\ell)} = F_k \setminus \{v_\ell\} \cup \{v_{a,1}, \ldots, v_{a,d}\}$$

where $d \geq 2$. First consider the sub-case where $\ell = r$. In this instance, $n \in \{i_r, j_r, k_r\}$. Without loss of generality we can say $n = i_r$. Since $B_r = \emptyset$, $i_{r-1} = n - 1$. Then we must have $n - 1 < i_{a,1} < i_{a,2}$ and $i_{a,1} = n$, but since $n$ is the maximum index allowed, this is a contradiction. Next, consider the sub-case where $\ell < r$. Then since $F_k$ satisfies P3, $\min\{i_{\ell+1} - i_\ell, j_{\ell+1} - j_\ell, k_{\ell+1} - k_\ell\} = 1$. Without loss of generality, we assume $i_{\ell+1} - i_\ell = 1$. Since $B_\ell = \emptyset$, $i_\ell - i_{\ell-1} = 1$. But we must have $i_{\ell-1} < i_{a,1} < i_{a,2} < i_\ell + 1$, which is impossible. Therefore we have also arrived at a contradiction when $\dim F_{j(\ell)} > \dim F_k$. So when $F_k$ is a homology facet, $B_\ell \neq \emptyset$ for all $\ell \in [r]$. □

Recall that for $n \geq 1$, $\Delta(n)$ has a collection of isolated vertices with index set of type $\{1, n, k\}$ where $k \in [n]$. Then, as a shellable simplicial complex has the homotopy type of a wedge of spheres, only a wedge including 0-spheres can increase the number of connected components to more than 1. In fact, we will obtain one additional connected component for each isolated vertex, so we refer to the isolated vertices of $\Delta(n)$ as *homology vertices*. Similarly, we say one-dimensional homology facets of $\Delta(n)$ are *homology edges*, and two-dimensional homology facets of $\Delta(n)$ are *homology triangles*. The next lemma shows how to count the homology vertices.

**Lemma 3.4.3.** *There are $6(n-1)$ homology vertices in the simplicial complex $\Delta(n)$.*

*Proof.* When $n = 2$, the 6 homology vertices are $(1,1,2), (1,2,1), (1,2,2), (2,1,1), (2,1,2)$, and $(2,2,1)$. For each $n$, there are 6 homology vertices of this form:

$$(1,1,n), (1,n,1), (1,n,n), (n,1,1), (n,1,n), (n,n,1).$$

Let $S_3$ denote the symmetric group on 3 letters. Then, for $n \geq 3$, for each $k \notin \{1, n\}$, there

are 6 homology vertices for the $|S_3| = 6$ vertices using the index set $\{1, k, n\}$. There are $n - 2$ choices for each such $k$, so there are $6 + 6(n - 2) = 6(n - 1)$ homology vertices in $\Delta(n)$. $\qquad\square$

### 3.4.1 Verifying the Identity 3.1 for $n \leq 4$

Recall that we hope to use the Euler-Poincaré relation to verify (3.1). We verify (3.1) for $n \leq 4$ by calculating the Betti numbers for $n \leq 4$. In these small cases, plugging $n$ into the relation

$$\sum_{i=-1}^{n-1} (-1)^{i+1} f_i = \sum_{i=-1}^{n-1} (-1)^{i+1} \beta_i \qquad (3.8)$$

which is Equation 3.5 multiplied on both sides by -1, shows the identity 3.1 holds.

The simplicial complex $\Delta(1)$ is a single point. So $\beta_{-1} = 1$, $\beta_0 = 1$ because there is a single connected component, and $\beta_k = 0$ for $k > 1$. Therefore the right hand side of (3.8) is $1 - 1 = 0$, and we have verified (3.1) for $n = 1$. For $\Delta(2)$, it is clear that we have 7 connected components (see Figure 3.1). So we have $\beta_{-1}(\Delta(2)) = 1$ and $\beta_0(\Delta(2)) = 7$. Then the right-hand side of Equation 3.8 is $1 - 7 = -6$. We evaluate (3.1) at $n = 2$:

$$(-1)^{2/2} \binom{3(2)/2}{(2)/2, (2)/2, (2)/2} = (-1)^2 \times 3! = -6.$$

We have verified (3.1) for $n = 2$.

We can use the homology vertices of $\Delta(2)$ to construct the homology edges of $\Delta(3)$, which by Lemma 3.4.2 are precisely the edges in $\Delta(3)$ satisfying $B_1 \neq \emptyset$ and $B_2 \neq \emptyset$. The edges of $\Delta(3)$ satisfying $B_1 \neq \emptyset$ are precisely the edges $\{v_1, v_2\}$ such that $v_1$ a homology vertex of $\Delta(2)$. Then for $v_1 \in \{(1, 2, 2), (2, 1, 2), (2, 2, 1)\}$ we must have $v_2 = (3, 3, 3)$ for $B_2 \neq \emptyset$ to hold. For the case $v_1 \in \{(1, 1, 2), (1, 2, 1)(2, 1, 1)\}$, we have 3 choices for each choice of $v_1$. We demonstrate for the 3 choices where $v_1 = (1, 1, 2)$: (1) $v_2 = (3, 3, 3)$, (2) $v_2 = (2, 3, 3)$, and (3) $v_3 = (3, 2, 3)$.

In particular, if $v_1$ has index set $\{1, 1, 2\}$, we obtain one homology edge by setting $v_2 = (3, 3, 3)$, and two homology edges for choosing a position $a_2 \in \{i_2, j_2, k_2\}$ to set equal to 2 for the corresponding position in $v_1 \setminus \{2\}$. So, there are 12 homology edges in $\Delta(3)$. By Lemma 3.4.3, there are 12 homology vertices, and therefore 13 connected components, in $\Delta(3)$, so $(-1) \times (-\beta_{-1} + \beta_0 - \beta_1) = 0$, verifying 3.1 for $n = 3$.

We can count the homology edges of $\Delta(4)$ similarly, using the homology vertices of $\Delta(2)$ and $\Delta(3)$ to construct homology edges of $\{v_1, v_2\} \in \Delta(4)$, as for both $n = 2$ and $n = 3$ the homology vertices give vertices $v_1 \in \Delta(4)$ satisfying $B_1 \neq \emptyset$. The next lemma facilitates these calculations.

**Lemma 3.4.4.** *Let $v = (i, j, k)$ be a homology vertex of $\Delta(3)$. The number of homology edges $\{v_1, v_2\}$ of $\Delta(4)$ such that $v = v_1$ is $(4 - i)(4 - j)(4 - k) - 1$.*

*Proof.* Since $(i, j, k)$ is a homology vertex of $\Delta(3)$, $3 \in \{i, j, k\}$ and the edge $F = \{(i, j, k), (4, 4, 4)\}$ satisfies P3 in Lemma 3.2.1. There are $4 - i$ ways to either make no change to $i_2$ or to replace $\{(i, j, k), (i_2, j_2, k_2)\}$ with the down-twist (Definition 3.2.4) $\{(i, j, k), (i_2 - 1, j_2, k_2)\}$. There are $4 - j$ ways to either make no change to $j_2$ or to replace $\{(i, j, k), (i_2, j_2, k_2)\}$ with the down-twist $\{(i, j, k), (i_2, j_2 - 1, k_2)\}$, and $4 - k$ ways to either make no change to $k_2$ or to replace $\{(i, j, k), (i_2, j_2, k_2)\}$ with the down-twist $\{(i, j, k), (i_2, j_2, k_2 - 1)\}$.

We obtain $(4 - i)(4 - j)(4 - k)$ edges this way. The down-twists will all be safe because $3 \in \{i, j, k\}$ so $4 \in \{i_2, j_2, k_2\}$ will hold. So, all $(4 - i)(4 - j)(4 - k)$ edges obtained thus far are facets. But the facet $\{(i, j, k), (i + 1, j + 1, k + 1)\}$, which is in the set of edges so obtained, does not satisfy $B_2 \neq \emptyset$ and so is not a homology edge by Lemma 3.4.2. Therefore there are $(4 - i)(4 - j)(4 - k) - 1$ homology edges $\{v_1, v_2\}$ of $\Delta(4)$ such that $v = v_1$. $\qquad\square$

We illustrate Lemma 3.4.4 in the next example.

**Example 3.4.5.** There are six homology vertices in $\Delta(3)$ for the six permutations of the index set $\{1, 2, 3\}$. For each of these vertices $v$ there are 5 homology edges $\{v_1, v_2\}$ such that $v = v_1$. Consider the case $v_1 = (1, 2, 3)$. We count the ways to obtain a homology edge by either setting $v_2 = (4, 4, 4)$ or obtaining a homology edge as a sequence of safe down-twists (Definition 3.2.4) of a homology edge. We have $4 - i_1 = 3$, $4 - j_1 = 2$, and $4 - k_1 = 1$. Then we can obtain three facets by either not changing $i_2$ or by reducing the index $i_2$ by 1 in a safe down-twist about $v_2$ in the $i$-column:

$$\{(1, 2, 3), (4, 4, 4)\}, \{(1, 2, 3), (3, 4, 4)\}, \{(1, 2, 3), (2, 4, 4)\}.$$

From these, we can obtain the following safe down-twists in the $j$-column:

$$\{(1, 2, 3), (4, 3, 4)\}, \{(1, 2, 3), (3, 3, 4)\}, \{(1, 2, 3), (2, 3, 4)\}.$$

Then from each of these we cannot make any moves in the $k$-column. Note that the facet $\{(1, 2, 3), (2, 3, 4)\}$ is not a homology facet because it does not satisfy $B_2 \neq \emptyset$. So there are $(4 - i_1)(4 - j_1)(4 - k_1) - 1 = (3)(2) - 1 = 5$ homology edges of $\Delta(4)$ such that $v_1 = (1, 2, 3)$.

Applying Lemma 3.4.4, we have 5 edges for each of the six homology vertices of $\Delta(3)$ with index set $\{1, 2, 3\}$, 8 homology edges for each of the 3 homology vertices of $\Delta(3)$ with index set $\{1, 1, 3\}$, and 2 homology edges for each of the 3 homology vertices of $\Delta(3)$ with index set $\{1, 3, 3\}$. So there are $5(6) + 8(3) + 2(3) = 60$ homology edges in $\Delta(4)$ with $v_1$ appearing as a homology vertex in $\Delta(3)$.

By inspection, we obtain 10 homology edges $\{v_1, v_2\}$ in $\Delta(4)$ where $v_1$ is one of the three homology vertices of $\Delta(2)$ with index set $\{1, 1, 2\}$, and we obtain 8 homology edges $\{v_1, v_2\}$ in

$\Delta(4)$ where $v_1$ is one of the three homology vertices of $\Delta(2)$ with index set $\{1, 2, 2\}$. So there are $10(3) + 8(3) = 54$ homology edges in $\Delta(4)$ with $v_1$ appearing as a homology vertex in $\Delta(2)$. In total, $\beta_1(\Delta(4)) = 114$.

The six homology triangles of $\Delta(4)$ are:

$$\{(1,1,2),(2,3,3),(4,4,4)\}, \{(1,1,2),(3,2,3),(4,4,4)\}, \{(1,2,1),(2,3,3),(4,4,4)\},$$

$$\{(1,2,1),(3,3,2),(4,4,4)\}, \{(2,1,1),(3,2,3),(4,4,4)\}, \{(2,1,1),(3,3,2),(4,4,4)\}.$$

Recall that by Lemma 3.4.3 there are $6(4-1) = 18$ homology vertices in $\Delta(4)$, so that $\beta_0(\Delta(4)) = 19$. Therefore the right hand side of (3.8) is $1 - 19 + 114 - 24 = 90$, which is the same as the right hand side of (3.1) evaluated at $n = 4$.

## 3.5   Discussion

We have not yet succeeded in counting the homology facets of $\Delta(n)$ for $n \geq 5$, either in individual cases or for general $n$. We know the number of homology vertices for all $n$ by Lemma 3.4.3. By Lemma 3.4.2, to count the $(r-1)$-dimensional homology facets for $r \geq 2$, it would be sufficient to count all facets $\{v_1, \ldots, v_r\}$ satisfying $B_\ell \neq \emptyset$ for all $\ell \in [r]$. However, this has not been accomplished at this point in time. Our inability to construct more general lemmas similar to Lemma 3.4.4 has precluded the development of a successful general approach counting homology facets of $\Delta(n)$ in terms of homology facets of $\Delta(k)$ for $k < n$.

As we noted in Section 3.1, the identity given in Equation 3.4 appears in Chapter 5 of Martin Aigner's book [1]. Chapter 5 of this book is titled "The Involution Principle." We give a standard definition of a well-known tool used in involution methods here, known as a sign-reversing involution. Additional background on sign-reversing involutions can also be found in Chapter 2 of Richard Stanley's book *Enumerative Combinatorics, Volume I* [66].

**Definition 3.5.1.** Let $X$ be a finite set of objects and suppose each element of $X$ has a *sign*, or in other words has been assigned either the value $+1$ or the value $-1$ via a function $s$, so that $s(x) \in \{1, -1\}$ for all $x \in X$. A *sign-reversing involution* on $X$ is a function $f$ such that $f$ is a bijection $f : X \to X$ satisfying

- $f(f(x)) = x$ for all $x \in X$, and

- if $f(x) \neq x$, then $s(f(x)) = -s(x)$.

Given a sign-reversing involution $f : X \to X$, let $X_F$ denote the set of all elements of $X$ such that $f(x) = x$. Then

$$\sum_{x \in X_F} s(x) = \sum_{x \in X} s(x).$$

Therefore, to total all signed elements of $X$ it is sufficient to find the cardinality of fixed points under $f$ which is the set $X_F$. Both (3.3) and (3.4) can be established using a sign-reversing involution. It is possible that we could complete our new proof of (3.1) using involution techniques on the homology facets of $\Delta(n)$. While we have tried to find a sign-reversing involution on the homology facets of $\Delta(n)$ towards the goal of reducing the total number of objects to be counted, we have not yet made progress with this strategy.

Furthermore, as we mentioned in Section 3.1, $\Delta(n)$ is the order complex of $P_{\Delta(n)}$ where the elements of $P_{\Delta(n)}$ are integer triples in $[n]^3$ and $(i, j, k) \leq (i', j', k')$ in $P_{\Delta(n)}$ if and only if $i < i'$, $j < j'$, and $k < k'$ as integers. This suggests another potential strategy towards computing the Betti numbers of $\Delta(n)$. Let $\widehat{P_{\Delta(n)}}$ be the poset obtained by adjoining a $\hat{0}$ and a $\hat{1}$ to $P_{\Delta(n)}$. If one were able to find a CL-labeling (Definition 2.3.1) of $\widehat{P_{\Delta(n)}}$, the CL-labeling could be used to compute the Betti numbers of $\Delta(n)$, as $\beta_r(\Delta(n))$ is the number of decreasing chains of length $r + 2$ in a CL-labeling of the Hasse diagram of $\widehat{P_{\Delta(n)}}$. See Section 5, and in particular Theorem 5.9 of the paper [15] by Anders Björner and Michelle Wachs for a detailed exploration of this approach.

Since EL-shellability implies CL-shellability, as we discussed in Section 2.3, an EL-labeling of $\widehat{P_{\Delta(n)}}$ would also be sufficient to compute the Betti numbers of $\Delta(n)$. Finally, an EL- or CL-labeling of either $\widehat{P_{\Delta(n)}}$ or $P_{\Delta(n)}$ could lead to a different shelling order for $\Delta(n)$ than the order $\mathcal{O}$ established in Theorem 3.3.4, which could in turn provide more insight into the structure of $\Delta(n)$.

# Chapter 4

# Polyhedral Combinatorics of UPGMA Cones

## 4.1 Introduction

In this chapter we study the decomposition of the input space $\mathbb{R}_{\geq 0}^{n(n-1)/2}$ induced by the distance-based phylogenetic reconstruction algorithm UPGMA. The UPGMA algorithm (Unweighted Pair Group Method with Arithmetic Mean) of Robert Sokal and Charles Michener (Algorithm 4.2.1), which was first proposed in the papers [60] and [61], is a clustering algorithm and distance-based phylogenetic tree reconstruction method. Like all distance-based phylogenetic reconstruction methods (Definition 1.4.16), UPGMA takes a dissimilarity map $\delta \in \mathbb{R}^{n(n-1)/2}$ (Definition 1.4.8) and returns a tree metric $d$ (Definition 1.4.9), which we can identify as a point in $\mathbb{R}^{n(n-1)/2}$. In practical use, it suffices to consider the action of UPGMA on inputs with nonnegative entries only, and when the entries of the input vector are nonnegative, the edge weights in the tree realization of the output $d$ will also be nonnegative. UPGMA outputs ultrametrics (Definition 1.4.13), which are equivalent to equidistant, rooted tree metrics (Definition 1.4.12).

Therefore we will study the partition

$$\{C(T) : \text{UPGMA}(x) \text{ is realized by } T \text{ for all } x \in C(T)\}$$

of $\mathbb{R}_{\geq 0}^{n(n-1)/2}$ induced by UPGMA, where $T$ is the binary rooted ranked tree that realizes the ultrametric output of UPGMA for all $\delta \in C(T)$. A motivation for studying the geometry of the regions in this partition is to gain insight into the observed performance of the algorithm on real datasets. Note that UPGMA, like its relative the Neighbor-Joining Algorithm (NJ) (Algorithm 5.4.6), can be viewed as a greedy heuristic for the NP-hard least-squares phylogeny problem

(LSP) (Problem 5.1.1). The relationship between the UPGMA and NJ algorithms and LSP is discussed in detail in Chapter 5.

Another motivation for studying the UPGMA algorithm is described in the paper [2], where David Aldous observed that rooted trees constructed from biological data using phylogenetic reconstruction methods are not distributed according to probability distributions similar to those followed by trees simulated under speciation models. As an example of a popular speciation model, we give the *Yule-Harding model*. The Yule-Harding model generates random rooted, binary phylogenetic trees with leaf set $[n]$ according to the following algorithm, commonly known as the *Yule process*. Recall that a *pendant edge* is an edge in a phylogenetic tree incident to a leaf.

**Algorithm 4.1.1** (Yule Process).    • Input: Leaf set $[n] = \{1, \ldots, n\}$, $n \geq 2$.

- Output: a rooted, binary phylogenetic $[n]$-tree.

- Initialize: Randomly select two leaves $x$ and $y$ from $[n]$ with uniform probability. Identify these two leaves as the leaf set of a rooted binary tree $T_1$. Set $S_1 = [n] \setminus \{x, y\}$ and $L_1 = \{x, y\}$.

- While $S_i \neq \emptyset$:

    - Randomly select an element $x_i$ of $S_i$ with uniform probability.
    - Randomly select a pendant edge $e_i = (u_i, y_i)$ of $T_i$ with probability determined by the uniform distribution on the set of pendant edges of $T_i$. Here $u_i$ is a binary internal vertex and $y_i \in L_i$.
    - Subdivide $e_i$ by adding a new vertex $v_i$.
    - Update $T_{i+1}$ as the tree with leaf set $L_{i+1} = L_i \cup \{x_i\}$ , $V(T_{i+1}) = V(T_i) \cup \{v_i, x_i\}$ and $E(T_{i+1}) = E(T_i) \setminus \{e_i\} \cup \{(u_i, v_i), (v_i, y_i), (v_i, x_i)\}$. Update $S_{i+1} = S_i \setminus \{x_i\}$.

- Return $T_{n-2+1}$, a rooted, binary phylogenetic tree with leaf set $[n]$.

Aldous observes in [2] that trees reconstructed from data are less balanced than trees predicted by the Yule-Harding model, where the balance of a tree is expressed in terms of the average number of descendants of each interior vertex. The rooted tree shape or tree topology (Definition 1.4.2) that is the least balanced is called the *comb* topology. Figure 4.1 shows a comb with six leaves. The relationship between the Yule-Harding model and other speciation models is discussed in Section 2.5 of *Phylogenetics* by Charles Semple and Mike Steel [57]. The lack of agreement in balance between simulated trees and trees reconstructed from biological data leads us to consider the notion of potential bias in reconstruction methods.
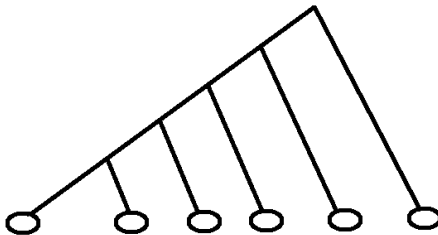
Figure 4.1: A comb tree with six leaves.

The goal of our analysis of the partition of $\mathbb{R}_{\geq 0}^{n(n-1)/2}$ induced by UPGMA is to gain insight into this problem of potential bias from the perspective of geometric combinatorics: if regions corresponding to some tree shapes are inherently larger than others, this indicates that UPGMA might favor those shapes in the presence of noise in data or model misspecification of the equidistant assumption (Definition 1.4.12) on the tree metric representing the true evolutionary history.

With these motivating problems in mind, we study the decomposition of space induced by the UPGMA algorithm. For a given binary phylogenetic $X$-tree $T$ (that is, with leaf labels $X$ but without edge weights), the region $\mathcal{P}(T) \subseteq \mathbb{R}_{\geq 0}^{n(n-1)/2}$ of dissimilarity maps for which the algorithm returns $T$ is a union of finitely many polyhedral cones, one for each rank function (Definition 1.4.4) on the interior nodes of $T$. We give explicit polyhedral descriptions of the cones, including $\mathcal{H}$-representations and $\mathcal{V}$-representations, for all $T$ and all $n$. In particular, each cone has $O(n^3)$ facet-defining inequalities in the $\mathcal{H}$-representation but exponentially many extreme rays in the $\mathcal{V}$-representation.

To compare the relative sizes of the regions $\mathcal{P}(T)$, we use the idea of spherical volume:

**Definition 4.1.2.** Let $C$ be a cone in $\mathbb{R}^d$. The *spherical volume* of $C$ is the surface area of the intersection of $C$ with the $(d-1)$-dimensional unit sphere $S_{d-1}$ in $\mathbb{R}^d$.

For example, if $C$ is the positive quadrant in $\mathbb{R}^2$, then the spherical volume of $C$ is equal to the angle measure $\pi/2$, which is the 1-dimensional area of the intersection of $S_1$ with $C$. We estimate the spherical volumes of the regions $\mathcal{P}(T)$ for $n \leq 7$. These volumes give a measure of the proportion of dissimilarity maps for which UPGMA returns a given combinatorial type of tree. In particular, our computations seem to indicate that highly unbalanced trees have small volume UPGMA cones compared to more balanced trees. Our computation of spherical volumes

builds on the Monte Carlo strategy in the paper [31] of Kord Eickmeyer, Peter Huggins, Lior Pachter, and Ruriko Yoshida. The results in this chapter are joint work with Seth Sullivant, and appear in the paper [24].

## 4.2 Ranked Phylogenetic Trees and the UPGMA Algorithm

Recall that the lattice of set partitions $\Pi_n$ provides a useful alternate description of rooted, ranked phylogenetic trees. As we saw in Example 1.4.7, every maximal chain in the lattice of set partitions corresponds to a ranked phylogenetic tree. For a maximal chain $C$ in $\Pi_n$ we write

$$C = 1|2|\cdots|n = \pi_n \lessdot \pi_{n-1} \lessdot \cdots \lessdot \pi_2 \lessdot \pi_1 = 12\cdots n.$$

We use the convention that $\pi_i$ is always a partition with $i$ parts.

Given $\pi_i \in C$, we write $\pi_i = \lambda_1^i|\lambda_2^i|\cdots|\lambda_i^i$. When $\pi_i \lessdot \pi_{i-1}$, there are exactly two blocks $\lambda_j^i, \lambda_k^i$ that are joined in $\pi_{i-1}$ but distinct in $\pi_i$. Recall that $\mathring{V}(T)$ denotes the set of interior vertices of $T$. If $v \in \mathring{V}(T)$ where $r(v) = n - i$, then $\pi_{i-1}$ joins the two blocks in $\pi_i$ that correspond to the subtrees of $T$ induced by the child nodes of $v$.

The UPGMA algorithm constructs a rooted ranked phylogenetic $X$ tree from a dissimilarity map $\delta$, as well as an equidistant tree metric $d$ which approximates $\delta$. The algorithm works as follows:

**Algorithm 4.2.1** (UPGMA Algorithm). • Input: a dissimilarity map $\delta \in \mathbb{R}_{\geq 0}^{n(n-1)/2}$ on $[n]$.

- Output: a maximal chain $C$ in the partition lattice $\Pi_n$ and an equidistant tree metric $d$.

- Initialize $\pi_n = 1|2|\cdots|n$, and set $\delta^n = \delta$.

- For $i = n - 1, \ldots, 1$ do

  - From partition $\pi_{i+1} = \lambda_1^{i+1}|\cdots|\lambda_{i+1}^{i+1}$ and distance vector $\delta^{i+1} \in \mathbb{R}_{\geq 0}^{(i+1)i/2}$ choose $j, k$ be so that $\delta^{i+1}(\lambda_j^{i+1}, \lambda_k^{i+1})$ is minimized.

  - Set $\pi_i$ to be the partition obtained from $\pi_{i+1}$ by merging $\lambda_j^{i+1}$ and $\lambda_k^{i+1}$ and leaving all other parts the same. Let $\lambda_i^i = \lambda_j^{i+1} \cup \lambda_k^{i+1}$.

  - Create a new dissimilarity map $\delta^i \in \mathbb{R}_{\geq 0}^{i(i-1)/2}$ by setting $\delta^i(\lambda, \lambda') = \delta^{i+1}(\lambda, \lambda')$ if $\lambda, \lambda'$ are both parts of $\pi_{i+1}$ and

$$\delta^i(\lambda, \lambda_i^i) = \frac{|\lambda_j^{i+1}|}{|\lambda_i^i|}\delta^{i+1}(\lambda, \lambda_j^{i+1}) + \frac{|\lambda_k^{i+1}|}{|\lambda_i^i|}\delta^{i+1}(\lambda, \lambda_k^{i+1})$$

  otherwise.

– For each $x \in \lambda_j^{i+1}$ and $y \in \lambda_k^{i+1}$, set $\delta^i(x,y) = \delta^{i+1}(\lambda_j^{i+1}, \lambda_k^{i+1})$.

• Return: Chain $C = \pi_n \lessdot \cdots \lessdot \pi_1$ and equidistant tree metric $d = \delta^2$.

Note that at the step which recalculates distances, the weighted average

$$\delta^i(\lambda, \lambda_i^i) = \frac{|\lambda_j^{i+1}|}{|\lambda_i^i|} \delta^{i+1}(\lambda, \lambda_j^{i+1}) + \frac{|\lambda_k^{i+1}|}{|\lambda_i^i|} \delta^{i+1}(\lambda, \lambda_k^{i+1})$$

is used to determine the new distance. This is simply a computationally efficient strategy to compute the average

$$\delta^i(\lambda, \lambda') = \frac{1}{|\lambda| \cdot |\lambda'|} \sum_{x \in \lambda, y \in \lambda'} \delta(x,y) \tag{4.1}$$

a formula we will make use of later.

**Example 4.2.2.** Let $\delta = (1, 2, 1.8, 1.7, 2, 2.6, 3.1, 2.4, 2.6, 1.2) \in \mathbb{R}_{\geq 0}^{5(5-1)/2}$ be a dissimilarity map on 5 taxa.

The UPGMA algorithm performs the following steps, where an underline is used to denote the smallest value in $\delta^i$.

$$\delta^5 = \begin{array}{cccccccccc} 12 & 13 & 14 & 15 & 23 & 24 & 25 & 34 & 35 & 45 \\ (\underline{1}, & 2, & 1.8, & 1.7, & 2, & 2.6, & 3.1, & 2.4, & 2.6, & 1.2) \end{array}$$

$$\delta^4 = \begin{array}{cccccc} 12,3 & 12,4 & 12,5 & 34 & 35 & 45 \\ (2, & 2.2, & 2.4, & 2.4, & 2.6, & \underline{1.2}) \end{array}$$

$$\delta^3 = \begin{array}{ccc} 12,3 & 12,45 & 3,45 \\ \underline{2} & 2.3 & 2.5 \end{array}$$

$$\delta^2 = \begin{array}{c} 123,45 \\ \underline{2.367} \end{array}$$

where

$$2.367 \approx \left( \frac{|12|}{|12| + |3|} \right)(2.3) + \left( \frac{|3|}{|12| + |3|} \right)(2.5)$$

The resulting ranked, rooted tree $T$ with an equidistant weighting $w$ produced by the UPGMA algorithm is displayed in Figure 4.2. Note that the rank function on the tree corresponds to the agglomeration steps in the algorithm. The weighting $w$ is equidistant because at step $i$ of the algorithm, the distance between every pair of elements in two blocks $\lambda, \lambda'$ in $\pi_i$ is the weighted average of all distances between all pairs, as shown in Equation 4.1.
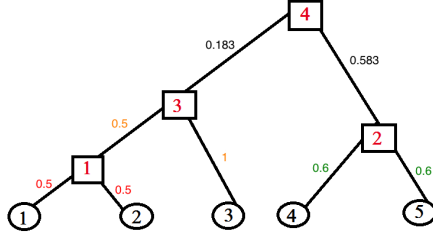
Figure 4.2: The tree metric $d$ realized by the weighted tree $(T, w)$ from Example 4.2.2.

The corresponding chain in the lattice of partitions $\Pi_5$ is

$$C = 1|2|3|4|5 \lessdot 3|4|5|12 \lessdot 3|12|45 \lessdot 45|123 \lessdot 12345.$$

## 4.3 UPGMA Regions and UPGMA Cones

The UPGMA algorithm partitions the set of dissimilarity maps with positive entries, which we identify as $\mathbb{R}_{\geq 0}^{n(n-1)/2}$, into regions indexed by ranked, rooted tree shapes corresponding to maximal chains in $\Pi_n$. For a given rooted phylogenetic $[n]$-tree $T$ let $\mathcal{P}(T) \subseteq \mathbb{R}_{\geq 0}^{n(n-1)/2}$ denote the Euclidean topological closure of the set of dissimilarity maps $\delta$ such that the UPGMA algorithm returns $T$ when given $\delta$ as an input. The set $\mathcal{P}(T)$ is called the *UPGMA region* associated to the tree $T$. Similarly, for a maximal chain $C$ in $\Pi_n$, let $\mathcal{P}(C) \subseteq \mathbb{R}_{\geq 0}^{n(n-1)/2}$ denote the closure of the set of dissimilarity maps such that the UPGMA algorithm returns the ranked rooted tree corresponding to the chain $C$.

Our goal in this section is to describe the sets $\mathcal{P}(T)$ and $\mathcal{P}(C)$. Clearly $\mathcal{P}(T) = \cup \mathcal{P}(C)$ where the union is over all maximal chains in $\Pi_n$ whose associated tree is $T$. The next theorem, Theorem 4.3.1, describes how each $\mathcal{P}(T)$ is a polyhedral cone, and the steps in the algorithm UPGMA determine the $\mathcal{H}$-representation of each cone $\mathcal{P}(C)$. Recall that a pointed cone has a trivial lineality space (Definition 1.2.7) consisting only of the origin.

**Theorem 4.3.1.** *For each chain $C \in \Pi_n$ the set $\mathcal{P}(C)$ is a pointed polyhedral cone with $O(n^3)$ facet-defining inequalities, and exponentially many extreme rays. Each covering relation in the chain $C$ determines a collection of facet-defining inequalities for $\mathcal{P}(C)$. Each element of the chain $C$ determines a collection of extreme rays of $\mathcal{P}(C)$.*

Theorem 4.3.1 is of practical interest concerning UPGMA, but we will provide a more general result for the description of cones associated to partial chains, which we define next, that will imply Theorem 4.3.1.

**Definition 4.3.2.** A *partial chain C* is a sequence

$$\pi_s \lessdot \pi_{s-1} \lessdot \cdots \lessdot \pi_t$$

for some $n \geq s \geq t \geq 1$. A partial chain is *grounded* if $s = n$.

The fact that steps in the partial chains correspond to covering relations guarantees that at each step the covering relation $\pi_{i+1} \lessdot \pi_i$ corresponds to joining a pair of blocks of the partition $\pi_{i+1}$ into a new block of $\pi_i$. This means that any partial chain $C$ can be interpreted as intermediate information that is calculated between steps $s$ and $t$ of the UPGMA algorithm.

For a partial chain $C$, let $\mathcal{P}(C)$ denote the set of all dissimilarity maps in $\mathbb{R}_{\geq 0}^{s(s-1)/2}$ which the UPGMA algorithm could produce on steps $s$ through $t$ of the algorithm. Observe that this notation is consistent with the use of $\mathcal{P}(C)$ to denote the cone corresponding to a maximal chain $C \in \Pi_n$. We let the coordinates in the space $\mathbb{R}^{s(s-1)/2}$ represent the $s(s-1)/2$ distances $\delta(\lambda_j^s, \lambda_k^s)$.

**Proposition 4.3.3.** *Let $C$ be a partial chain in $\Pi_n$. Let $\mathcal{P}(C) \subseteq \mathbb{R}_{\geq 0}^{s(s-1)/2}$ be the set of dissimilarity maps for which steps $s$ through $t$ of the UPGMA algorithm return the partial chain $C$. For each covering relation $\pi_i \lessdot \pi_{i-1}$, let $\lambda_{j(i)}^i$ and $\lambda_{k(i)}^i$ be the pair of blocks of $\pi_i$ joined in $\pi_{i-1}$. Then $\mathcal{P}(C)$ is the solution to the following system of linear inequalities:*

$$\delta(\lambda_j^s, \lambda_k^s) \geq 0 \text{ for all } j, k$$

$$\text{for } i = s, \ldots, t-1, \text{ and for all pairs } j, k \neq j(i), k(i)$$

$$\frac{1}{|\lambda_{j(i)}^i||\lambda_{k(i)}^i|} \sum_{\lambda_j^s \subseteq \lambda_{j(i)}^i, \lambda_k^s \subseteq \lambda_{k(i)}^i} |\lambda_j^s||\lambda_k^s|\delta(\lambda_j^s, \lambda_k^s) \leq \frac{1}{|\lambda_j^i||\lambda_k^i|} \sum_{\lambda_j^s \subseteq \lambda_j^i, \lambda_k^s \subseteq \lambda_k^i} |\lambda_j^s||\lambda_k^s|\delta(\lambda_j^s, \lambda_k^s).$$

*Note that if $s > t$ we only need the nonnegativity constraint $\delta(\lambda_{j(s)}^s, \lambda_{k(s)}^s) \geq 0$, as the other inequalities $\delta(\lambda_j^s, \lambda_k^s) \geq 0$ follow from $\delta(\lambda_{j(s)}^s, \lambda_{k(s)}^s) \leq \delta(\lambda_j^s, \lambda_k^s)$.*

*Proof.* At step $i$ of the UPGMA algorithm, we choose the pair of $\lambda_{j(i)}^i$ and $\lambda_{k(i)}^i$ to merge such that $\delta^i(\lambda_{j(i)}^i, \lambda_{k(i)}^i)$ is minimized. Using the formula from Equation 4.1

$$\delta^i(\lambda_j^i, \lambda_k^i) = \frac{1}{|\lambda_j^i||\lambda_k^i|} \sum_{x \in \lambda_j^i, y \in \lambda_k^i} \delta(x, y)$$

twice shows that

$$\delta^i(\lambda_j^i, \lambda_k^i) = \frac{1}{|\lambda_j^i||\lambda_k^i|} \sum_{\lambda_j^s \subseteq \lambda_j^i, \lambda_k^s \subseteq \lambda_k^i} |\lambda_j^s||\lambda_k^s|\delta(\lambda_j^s, \lambda_k^s).$$

This yields precisely the inequalities in the statement of the proposition at step $i$. □

**Proposition 4.3.4.** *Given a maximal chain $C \in \Pi_n$, there are $O(n^3)$ facet-defining inequalities in the $\mathcal{H}$-representation of the polyhedral cone $\mathcal{P}(C)$.*

*Proof.* At step $t$, there are $\binom{t}{2}$ ways to merge two blocks of $\pi_t$, and the pair of blocks $\delta(\lambda_{j(t)}^t, \lambda_{k(t)}^t)$ merged at step $t$ can be paired with $\binom{t}{2} - 1$ other pairs of blocks. So $\binom{t}{2} - 1$ new inequalities are introduced at step $t$. An elementary identity for binomial coefficients tells us that for $a, b \geq 0$,

$$\sum_{r=b}^{a} \binom{r}{b} = \binom{a+1}{b+1}.$$

Thus there are

$$\sum_{t=2}^{n} \left( \binom{t}{2} - 1 \right) = \binom{n+1}{3} - n + 1$$

facet-defining inequalities. $\qquad\qquad\square$

According to Theorem 1.2.5, the facet-defining inequalities of the polyhedral cones $\mathcal{P}(C)$ imply the existence of a closed-form description of the extreme rays of $\mathcal{P}(C)$. Therefore, we now provide a description of the extreme rays of the cones of grounded partial chains $\mathcal{P}(C)$, which are partial chains (Definition 4.3.2) starting with the bottom element $\pi_n = 1|2|\cdots|n$. The polyhedral description of the cones $\mathcal{P}(C)$ for more general partial chains is used in the proof of Theorem 4.3.7, which gives the combinatorial description of the $\mathcal{V}$-representation of the extreme rays of $\mathcal{P}(C)$. We require some additional terminology to state the relevant results, beginning with Definition 4.3.5.

**Definition 4.3.5.** Given a partition $\pi_k = \lambda_1|\lambda_2|\cdots|\lambda_k \in \Pi_n$ a *traversal* of $\pi_k$ is a subset $F \subset \binom{[n]}{2}$ of size $\binom{k}{2}$, where each element of $F$ is a pair $\{p, p'\} \in \pi$ satisfying $p \in \lambda, p' \in \lambda'$. In each traversal there is precisely one such pair $p, p'$ for every pair of distinct blocks $\lambda, \lambda'$ of $\pi_k$.

For example, the partition $12|3|45$ has $2^2 \cdot (2 \cdot 1) \cdot (2 \cdot 1) = 16$ traversals. Figure 4.3 shows the traversal $F = \{\{1,3\}, \{1,4\}, \{3,5\}\}$ of the partition $12|3|45$, where for example the pairing $\{1,3\}$ is indicated by connecting the circled elements 1 and 3 with an edge. The purpose of introducing the language of traversals is to give a convenient method of describing the extreme rays, or $\mathcal{V}$-representation of $\mathcal{P}(C)$, for arbitrary chains $C \subset \Pi_n$.

**Definition 4.3.6.** Let $\pi_k = \lambda_1|\lambda_2|\cdots|\lambda_k \in \Pi_n$. Let $F$ be a traversal of $\pi_k$. The *induced vector* of $F$, denoted $v(F)$, is the dissimilarity map in $\mathbb{R}^{n(n-1)/2}$ such that

1. $v(F)(i,j) = 0$ if the pair $i, j$ is not in the traversal $F$, and

2. when $\{i, j\} \in F$, $v(F)(i,j) = |\lambda_{k(i)}||\lambda_{k(j)}|$ where $i \in \lambda_{k(i)}$ and $j \in \lambda_{k(j)}$.
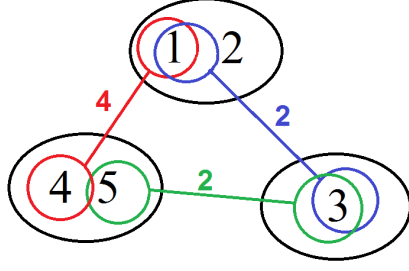
Figure 4.3: The traversal $F = \{\{1,3\},\{1,4\},\{3,5\}\}$ of the partition $12|3|45$.

This traversal $F = \{\{1,3\},\{1,4\},\{3,5\}\}$ of the partition $12|3|45$ shown in Figure 4.3 induces the vector $(0,2,4,0,0,0,0,0,2,0)$. For $i \neq j \in [n]$ let $e_{ij} \in \mathbb{R}^{n(n-1)/2}$ be the dissimilarity map such that $e_{ij}(i,j) = 1$ and $e_{ij}(x,y) = 0$ for all other pairs $x \neq y$.. Another way to write $v(F)$ is

$$v(F) = 2 \cdot e_{1,3} + 4 \cdot e_{1,4} + 2 \cdot e_{3,5}.$$

**Theorem 4.3.7.** *Let $C = \pi_n \lessdot \pi_{n-1} \lessdot \cdots \lessdot \pi_t$ be a grounded partial chain in $\Pi_n$. Then $\mathcal{P}(C)$ is a cone with extreme rays given by the set of vectors*

$$\{e_{k,l} : k, l \text{ are not in the same part of the partition } \pi_t\}$$
$$\bigcup \; \bigcup_{i=t+1}^{n} \{v(F) \; : \; F \text{ is a traversal of } \pi_i \}.$$

Note that if $t = 1$, the first set in the union is empty.

The remainder of this section consists of the proof of Theorem 4.3.7 and completes our description of the cones $\mathcal{P}(C)$. The proof will be broken into a number of pieces, and will work by induction on both $t$ and $n$. Let $\mathbf{1}^t$ denote the vector in $\mathbb{R}^{t(t-1)/2}$ all of whose coordinates are equal to one. Note that $\mathbf{1}^n$ is the induced vector of the single traversal associated to the partition $1|2|\cdots|n$, which appears in every partial chain.

**Lemma 4.3.8.** *Let $C = \pi_s \lessdot \cdots \lessdot \pi_t$ be a partial chain in $\Pi_n$ with $s > t$. Then*

1. *$\mathbf{1}^s$ is an extreme ray of $\mathcal{P}(C)$, and*

2. *$\mathbf{1}^s$ is the only extreme ray of $\mathcal{P}(C)$ that has a nonzero $(\lambda^s_{j(s)}, \lambda^s_{k(s)})$ coordinate where $(\lambda^s_{j(s)}, \lambda^s_{k(s)})$ is the pair of parts joined together in the partition $\pi_{s-1}$.*

*Proof.* First of all, $\mathbf{1}^s$ satisfies the single inequality $\delta(\lambda_{j(s)}, \lambda_{k(s)}) \geq 0$ of Proposition 4.3.3 strictly. The rest of the inequalities of Proposition 4.3.3 are satisfied with equality by $\mathbf{1}^s$ so that $\mathbf{1}^s \in \mathcal{P}(C)$. Hence the extreme ray $\mathbf{1}^s$ is in the intersection of all the facet-defining inequalities

except for one. Since $\mathcal{P}(C)$ is a pointed cone because it is contained in the positive orthant, this implies that $\mathbf{1}^s$ is an extreme ray. This proves part (1). Furthermore, since every extreme ray of a cone is the intersection of some of its facet-defining inequalities, every other extreme ray must have the inequality $\delta(\lambda_{j(s)}, \lambda_{k(s)}) \geq 0$ as an active inequality, meaning the inequality is actually an equality at every point on the extreme ray. This proves part (2). $\qquad\square$

Let $C = \pi_s \lessdot \cdots \lessdot \pi_t$ be a partial chain, and $C'$ a partial chain obtained as a *final segment* of $C$, that is, there is a $s < u \leq t$, such that $C' = \pi_u \lessdot \cdots \lessdot \pi_t$. The UPGMA algorithm induces a natural linear map $A(C, C') : \mathbb{R}^{s(s-1)/2} \to \mathbb{R}^{u(u-1)/2}$. In particular, it is defined by

$$(A(C, C')\delta)(\lambda, \lambda') = \frac{1}{|\lambda||\lambda'|} \sum_{\substack{\mu, \mu' \in \pi_s \\ \mu \subseteq \lambda, \mu' \subseteq \lambda'}} |\mu||\mu'|\delta(\mu, \mu')$$

where $\lambda, \lambda'$ are blocks of $\pi_u$.

**Definition 4.3.9.** A linear transformation $\phi : \mathbb{R}^n \to \mathbb{R}^m$ is a *coordinate substitution* if for each of the coordinate vectors $e_i$, $\phi(e_i) = c_i e_{\alpha(i)}$ with $c_i > 0$, where $\alpha : [n] \to [m]$. That is, each coordinate maps to a scaled version of another coordinate.

Note, in particular, the quantity $\delta(\mu, \mu')$ only appears in the formula for $(A(C, C')\delta)(\lambda, \lambda')$, so that $A(C, C')$ is a coordinate substitution map. when restricted to the coordinates $\delta(\mu, \mu')$ where $\mu, \mu'$ are in different parts of $\pi_s$.

With the preceding paragraph in mind, we let $\tilde{\mathcal{P}}(C)$ denote the intersection of $P(C)$ with the hyperplane $\{\delta : \delta(\lambda_{j(s)}, \lambda_{k(s)}) = 0\}$.

**Proposition 4.3.10.** *Let $C = \pi_s \lessdot \cdots \lessdot \pi_t$ be a partial chain and with final segment $C' = \pi_{s-1} \lessdot \cdots \lessdot \pi_t$. Then $A(C, C') : \tilde{\mathcal{P}}(C) \to \mathcal{P}(C')$ is surjective, and $\tilde{\mathcal{P}}(C) = A(C, C')^{-1}(\mathcal{P}(C')) \cap \mathbb{R}_{\geq 0}^{s(s-1)/2-1}$.*

*Proof.* Note that by definition of the UPGMA algorithm, the map $A(C, C') : \mathcal{P}(C) \to \mathcal{P}(C')$ is surjective. If a vector $\delta^s \in \mathcal{P}(C)$, then so is the vector

$$\delta' = \delta^s - \delta^s(\lambda_{j(s)}^s, \lambda_{k(s)}^s)e(\lambda_{j(s)}^s, \lambda_{k(s)}^s),$$

obtained by zeroing out the $(\lambda_{j(s)}^s, \lambda_{k(s)}^s)$ coordinate. However, $A(C, C')\delta^s = A(C, C')\delta'$, which implies that $A(C, C') : \tilde{\mathcal{P}}(C) \to \mathcal{P}(C')$ is surjective.

To see that $\tilde{\mathcal{P}}(C) = A(C, C')^{-1}(\mathcal{P}(C')) \cap \mathbb{R}_{\geq 0}^{s(s-1)/2-1}$, note that the inequalities that describe $\tilde{\mathcal{P}}(C)$ are precisely the pullbacks of the inequalities that describe $\mathcal{P}(C')$, plus nonnegativity constraints, since none of the inequalities on $\mathcal{P}(C)$ coming from the covering relation $\pi_s \lessdot \pi_{s-1}$ are needed. $\qquad\square$

**Lemma 4.3.11.** *Let $D \subseteq \mathbb{R}^m$ be a polyhedral cone, $\phi : \mathbb{R}^n \to \mathbb{R}^m$ be a coordinate substitution with associated map $\alpha$, and $C \subseteq \mathbb{R}^n$ a polyhedral cone such that $\phi(C) = D$. Suppose that $C = \mathbb{R}^n_{\geq 0} \cap \phi^{-1}(D)$. Let $V$ be the set of extreme rays of $C$. Then extreme rays of $D$ consist of all vectors obtained by the following procedure:*

*For each extreme ray $\sum_j a_j e_j \in V$, include all vectors of the form $\sum_j a_j/c_{\beta(j)} e_{\beta(j)}$ ranging over all functions $\beta : [m] \to [n]$ such that $\alpha(\beta(j)) = j$ for all $j$.*

*Proof.* It suffices to show that under the hypotheses of the Lemma, every extreme ray of $C$ maps onto an extreme ray of $D$. Indeed, if that is the case, the extreme rays of $C$ are precisely the vertices of the polytopes $\phi^{-1}(v) \cap \mathbb{R}^n_{\geq 0}$ as $v$ ranges over the extreme rays of $V$. Note that since $\phi$ is a coordinate substitution $\phi^{-1}(v)$ is isomorphic to a product of simplices, the simplices being defined over coordinate subsets over the form $\alpha^{-1}(j)$. The vertices of these products of simplices have the form of the statement of the Lemma.

Hence, it suffices to verify the claim that every extreme ray of $C$ maps onto an extreme ray of $D$. So suppose that $v'$ is an extreme ray of $C$ such that $\phi(v') = v$ is not an extreme ray of $D$. Then there exists $w, u \in D$, not equal to $v$ such that $v = w + u$. Using these vectors, we construct $w', u' \in C$ not equal to $v'$ such that $v' = w' + u'$. For each $i$ such that $\alpha(i) = j$ define

$$w'_i = \frac{w_j}{v_j} v'_i \qquad \text{and} \qquad u'_i = \frac{u_j}{v_j} v'_i.$$

Clearly with this choice, we have $v' = w' + u'$ since $v_j = w_j + u_j$, and both $w'$ and $u'$ consist of nonnegative vectors. Also, since $w, u$ are not equal to $v$, neither are $w', u'$ equal to $v'$. So we must show that $\phi(w') = w$ and $\phi(u') = u$. But

$$\phi(w')_j = \sum_{i:\alpha(i)=j} \frac{w_j}{v_j} c_i = \frac{w_j}{v_j} \sum_{i:\alpha(i)=j} c_i = \frac{w_j}{v_j} v_j = w_j.$$

A similar statement holds for $u'$, which completes the proof. $\square$

We now have all the ingredients to prove Theorem 4.3.7.

*Proof of Theorem 4.3.7.* Let $C = \pi_s \lessdot \cdots \lessdot \pi_t$. First of all, note that if $s = t$, then $\mathcal{P}(C)$ is the positive orthant in $\mathbb{R}^{s(s-1)/2}$, whose extreme rays are the standard unit vectors.

Now assume that $s > t$. According to Lemma 4.3.8, the vector $\mathbf{1}^s$ is an extreme ray of $\mathcal{P}(C)$. Letting $C' = \pi_{s-1} \lessdot \cdots \lessdot \pi_t$, Proposition 4.3.10 we see that all other extreme rays of $\mathcal{P}(C)$ can be obtained by applying Lemma 4.3.11 to the extreme rays of $\mathcal{P}(C')$. Repeating this procedure for the extreme rays of $\mathcal{P}(C)$ that do not map to $\mathbf{1}^{s-1} \in \mathcal{P}(C')$, we see that every extreme ray of $\mathcal{P}(C)$ besides $\mathbf{1}^s$ can either be obtained as a vertex of the polyhedron $A(C, C_u)^{-1}(\mathbf{1}^u)$ where $C_u = \pi_u \lessdot \cdots \lessdot \pi_t$, or as a vertex of $A(C, C_t)^{-1}(e_{\lambda_k, \lambda_l})$.

To complete the proof of the theorem we must analyze the vertices of $A(C, C_t)^{-1}(e_{\lambda_k, \lambda_l})$ and show that the vertices of $A(C, C_u)^{-1}(\mathbf{1}^u)$ are precisely the induced vectors from the traversals of $\pi_u$. For both of these statements, we can use Lemma 4.3.11.

Indeed, $A(C, C_u)$ is the map such that

$$(A(C, C_u)\delta)(\lambda, \lambda') = \frac{1}{|\lambda| \cdot |\lambda'|} \sum_{\substack{x \in \lambda \\ y \in \lambda'}} \delta(x, y).$$

This implies, by Lemma 4.3.11 that the vertices of

$$A(C, C_t)^{-1}(e_{\lambda, \lambda'})$$

are $|\lambda| \cdot |\lambda'| e_{k,l}$ such that $k \in \lambda$ and $l \in \lambda'$. Since we can ignore the scaling factor $|\lambda| \cdot |\lambda'|$ when describing extreme rays, taking the union over all pairs $\lambda, \lambda' \in \pi_t$, yields the set of rays $\{e_{k,l} : k, l \text{ are not in the same part of the partition } \pi_t\}$ from Theorem 4.3.7.

Similarly, applying Lemma 4.3.11 to the map $A(C, C_u)$ and the vector $\mathbf{1}^u$ yields the set of induced vectors associated to the partition $\pi_u$. Indeed, the coordinate 1 in the $(\lambda, \lambda')$ position of $\mathbf{1}^u$ produces an entry of $|\lambda| \cdot |\lambda'|$ in exactly one of the positions $\delta(x, y)$ such that $x \in \lambda, y \in \lambda'$. This completes the proof of Theorem 4.3.7. □

We now show that Theorem 4.3.7 implies that the UPGMA cones have exponentially many extreme rays.

**Proposition 4.3.12.** *The cones $\mathcal{P}(C)$ have exponentially many extreme rays.*

*Proof.* Given $\pi_s = \lambda_1^s | \cdots | \lambda_s^s$, the number of traversals is the product of the pairwise products of the cardinalities of the blocks of $\pi_s$. So the number of extreme rays induced by $\pi_s$ is

$$\prod_{\{i,j\} \subset \binom{[s]}{2}} |\lambda_i^s| |\lambda_j^s| = \prod_{i=1}^{s} |\lambda_i^s|^{s-1}.$$

Given a maximal chain $C \in \Pi_n$, the total number of extreme rays will be

$$\sum_{s=2}^{n} \prod_{i=1}^{s} |\lambda_i^s|^{s-1}$$

which is exponential in $n$. □

To summarize, we note that Propositions 4.3.3, 4.3.4, 4.3.12 and Theorem 4.3.7 yield Theorem 4.3.1.

## 4.4 Applications of Theorem 4.3.7

We use the characterization of the extreme rays in Theorem 4.3.7 of the cones $\mathcal{P}(C)$ to provide easy geometric applications. First, the set $\mathcal{P}(T)$ of all dissimilarity maps for which UPGMA returns a given tree is not a convex set in general. Second, the partition of the positive orthant into the cones $\mathcal{P}(C)$ does not have the structure of a polyhedral fan (Definition 1.2.8). Third, we show that among all possible tree topologies, the comb tree topology minimizes the number of rays in a UPGMA cone. A comb tree is pictured in Figure 4.1.

**Corollary 4.4.1.** *The UPGMA regions $\mathcal{P}(T) = \cup \mathcal{P}(C)$ are not convex in general.*

*Proof.* We give an example for $n = 4$. Let $T$ be the fork tree on 4 leaves shown in Figure 4.4. Then $\mathcal{P}(T) = \mathcal{P}(C_1) \cup \mathcal{P}(C_2)$ where

$$C_1 = 1|2|3|4 \lessdot 3|4|12 \lessdot 12|34 \lessdot 1234,$$

and

$$C_2 = 1|2|3|4 \lessdot 1|2|34 \lessdot 34|12 \lessdot 1234.$$

Now $v_1 = (0, 0, 2, 2, 0, 1)$ is an extreme ray of $P(C_1)$ induced by a traversal of $3|4|12$ and $v_2 = (1, 0, 2, 2, 0, 0)$ is an extreme ray of $\mathcal{P}(C_2)$ induced by a traversal of $1|2|34$. Let $\delta$ be the convex combination

$$\delta = \frac{1}{2}v_1 + \frac{1}{2}v_2 = \left(\frac{1}{2}, 0, 2, 2, 0, \frac{1}{2}\right).$$

If $\delta$ is input into UPGMA, the algorithm will return a tree with either $\{1, 3\}$ or $\{2, 4\}$ as the leaf set of a clade (Definition 1.4.5), so $\delta$ is not in $\mathcal{P}(T)$. So, in general, UPGMA regions are not convex unless $\mathcal{P}(T) = \mathcal{P}(C)$ for a single chain $C$ in $\Pi_n$. $\quad\square$
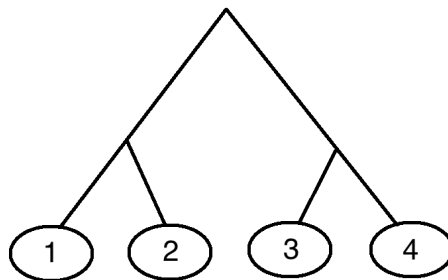


Figure 4.4: A fork tree with leaf set [4].

**Corollary 4.4.2.** *The UPGMA cones do not partition $\mathbb{R}^{n(n-1)/2}_{\geq 0}$ into a fan.*

*Proof.* Consider the two chains in $\Pi_4$

$$C_1 = \quad 1|2|3|4 \lessdot 3|4|12 \lessdot 4|123 \lessdot 1234$$

and

$$C_2 = \quad 1|2|3|4 \lessdot 2|4|13 \lessdot 4|123 \lessdot 1234.$$

The vector $(0, 0, 0, 1, 1, 1)$ generates an extreme ray of $P(C_1) \cap P(C_2)$ which we verified using the software polymake [37]. If $P(C_1) \cap P(C_2)$ were a face of $P(C_1)$ and $P(C_2)$, then $(0, 0, 0, 1, 1, 1)$ would generate a ray of $P(C_1)$ and $P(C_2)$. However, by Theorem 4.3.7, extreme rays of $P(C_1)$ and $P(C_2)$ must correspond to partitions in $\Pi_4$. Only partitions with 3 blocks induce vectors with 3 nonzero coordinates, and no partition of the set $[4]$ has 3 blocks of equal cardinality. So, no traversal of a partition in $\Pi_4$ induces a multiple of $(0, 0, 0, 1, 1, 1)$, and $P(C_1) \cap P(C_2)$ is not a face of a UPGMA cone. Recall from Definition 1.2.8 that two cones in a fan must intersect in a face of both. Therefore the UPGMA cones are not a fan. $\square$

The next corollary concerns the comb tree topology, which is the most unbalanced tree topology. Recall that Figure 4.1 shows a comb tree.

**Corollary 4.4.3.** *For each $n$, the comb tree topology minimizes the number of extreme rays over all UPGMA cones in $\mathbb{R}^{n(n-1)/2}_{\geq 0}$.*

*Proof.* Fix $n$. We will show that for each $1 \leq s \leq n$, the partitions whose parts have cardinalities $1, 1, \ldots, 1, n - s + 1$ minimize the number of traversals for all partitions with $s$ blocks. For all integers $x, y > 0$, we have $xy \geq (x + y - 1)(1)$. So for $\pi_s = \lambda_1^s | \cdots | \lambda_s^s$, the number of extreme rays induced by $\pi_s$ satisfies

$$\prod_{\{i,j\} \subset \binom{[s]}{2}} |\lambda_i^s| |\lambda_j^s| \geq \prod_{\{i,j\} \subset \binom{[s]}{2}} (1)(|\lambda_i^s| + |\lambda_j^s| - 1).$$

The only type of partition in $\Pi_n$ with $s$ parts such that all pairs $\{i, j\} \subset \binom{[s]}{2}$ satisfy either $|\lambda_s^i| = 1$ or $|\lambda_s^j| = 1$ is the type with $s - 1$ singleton blocks and one block of size $n - s + 1$. Therefore partitions of this type minimize the number of associated induced vectors.

If $C$ is a maximal chain in $\Pi_n$ such that every $\pi_s$ in $C$ is of this type, then the tree returned by $\delta \in \mathcal{P}(C)$ has the comb tree topology. Therefore this tree topology minimizes the possible number of extreme rays for any cone $\mathcal{P}(C) \subset \mathbb{R}^{n(n-1)/2}_{\geq 0}$. $\square$

## 4.5 Spherical Volumes of UPGMA Regions

A natural way to measure the region $\mathcal{P}(T)$ is to estimate the spherical volume (Definition 4.1.2) of the regions $\mathcal{P}(T) = \cup \mathcal{P}(C)$. In other words, we estimate the $\left(\binom{n}{2} - 1\right)$-dimensional area of the surface arising as the intersection of the cones $\mathcal{P}(C) \subset \mathcal{P}(T)$ with the unit sphere $S_{n(n-1)/2-1}$ in $\mathbb{R}_{\geq 0}^{n(n-1)/2}$.

We estimated the spherical volume of UPGMA cones in two ways using Mathematica, polymake [37], and the software [45] that accompanied the paper [31] of Kord Eickmeyer, Peter Huggins, Lior Pachter, and Ruriko Yoshida. For the first method, we sampled points from the positive orthant using a uniform spherical distribution and input the samples into UPGMA, recording which ranked tree the algorithm returned on the input point. The volume of $\mathcal{P}(T)$ is then the fraction of the total sample points returning $T$. We calculated volumes for $n = 4, 5, 6, 7$ using this method.

For the second method, we used a Monte Carlo integration strategy to estimate the surface area of the cones. The basic strategy can be described as follows. Recall that a cone in $\mathbb{R}^d$ is *simplicial* if the extreme rays of the cone are a linearly independent set of vectors in $\mathbb{R}^d$. Given a simplicial cone cone$(V)$ spanned by vectors $V = \{v_1, \ldots, v_k\}$, it is easy to generate uniform samples from the simplex conv$(V)$. The map that takes a point $x \in$ conv$(V)$ onto the surface cone$(V) \cap S_{n(n-1)/2-1}$ is simply $x \to x/\|x\|_2$. The spherical volume is then the average value of the Jacobian of this map. For $n = 4, 5, 6$, we used the software [45] to accomplish this method. This software requires as input triangulations of point configurations that we computed using the software polymake [37].

For $n = 7$, some triangulations for maximal chains in $\Pi_7$ were too large to compute and use. When it was not possible to compute the triangulation for the cone of a full chain, we computed a triangulation of a cone associated to a grounded partial chain $C'$ (Definition 4.3.2), used that triangulation to generate random samples, and then applied the UPGMA algorithm to see which chain was produced. In other words, we generated random points from the partial cone $\mathcal{P}(C')$ and computed the average of the product of the Jacobian and the indicator function of whether a point was in the cone $\mathcal{P}(C)$. We used Mathematica to implement this modification of the sampling strategy employed in [31].

Tables 4.1, 4.2, 4.3, and 4.4 summarize the results of our computations for $n = 4, 5, 6, 7$ leaf trees by displaying results for the regions $\mathcal{P}(T)$. So, these tables give estimates of the spherical volumes of the regions $\mathcal{P}(T)$. The column labeled Tree gives a representative of the combinatorial tree type in Newick format, which represents the clustering structure of the tree using parentheses. For example, the Newick format of the tree in Figure 4.2 is $(((12)3)(45))$, whereas $((((12)3)4)5)$ is a comb tree with leaf set [5]. The column labeled #Chains refers to the number of cones producing a fixed tree of the combinatorial type represented in the column

Table 4.1: Spherical volumes of UPGMA regions $\mathcal{P}(T)$ for $n = 4$.

| | Tree | # Chains | Volume | Fraction of Orthant |
|---|---|---|---|---|
| 1 | $(((12)3)4)$ | 1 | 0.0238 | 0.5895 |
| 2 | $((12)(34))$ | 2 | 0.0662 | 0.4099 |

labeled Tree. The column labeled Volume gives the total volume of all of the cones associated to the given tree, and the column labeled Fraction of Orthant gives the portion of the positive orthant in $\mathbb{R}^{n(n-1)/2}$ that returns the given tree type under UPGMA.

Recall that $\mathcal{P}(T) = \cup \mathcal{P}(C)$ where $C$ ranges over the chains in $\Pi_n$ corresponding to $T$. So, the number of cones associated to a tree $T$ depends on the number of rank functions that $T$ admits. For example, in the table for $n = 5$, the tree $T_2 = (((12)3)(45))$ has $4!/(4 \cdot 2 \cdot 1 \cdot 1) = 3$ rank functions by Equation 1.2, and so there are 3 cones in $\mathcal{P}(T_2)$.

The Mathematica software and input files for these computations are available at [25].

Table 4.2: Spherical volumes of UPGMA regions $\mathcal{P}(T)$ for $n = 5$.

| | Tree | # Chains | Volume | Fraction of Orthant |
|---|---|---|---|---|
| 1 | $((((12)3)4)5)$ | 1 | $8.57 \times 10^{-5}$ | 0.206 |
| 2 | $(((12)3)(45))$ | 3 | $5.01 \times 10^{-4}$ | 0.604 |
| 3 | $(((12)(34))5)$ | 2 | $3.14 \times 10^{-4}$ | 0.189 |

## 4.6  Discussion

The spherical volume computations suggest some conjectures which might hold true for large $n$. As Corollary 4.4.3 shows, the cone associated to the single rank function on the comb tree yields the cone $\mathcal{P}(C)$ with the fewest number of extreme rays. Our computations up to $n = 7$ suggest that this is also the cone with the smallest spherical volume. The size of the region $\mathcal{P}(T)$ appears to be roughly proportional to the number of chains $C$ that yield the tree $T$ and appears to be smallest for the comb tree. Furthermore, the relative proportion of the positive

Table 4.3: Spherical volumes of UPGMA regions $\mathcal{P}(T)$ for $n = 6$.

|   | Tree | # Chains | Volume | Fraction of Orthant |
|---|---|---|---|---|
| 1 | $(((((12)3)4)5)6)$ | 1 | $2.05 \times 10^{-8}$ | 0.042 |
| 2 | $((((12)3)4)(56))$ | 4 | $2.10 \times 10^{-7}$ | 0.216 |
| 3 | $((((12)3)(45))6)$ | 3 | $2.16 \times 10^{-7}$ | 0.223 |
| 4 | $(((12)3)((45)6))$ | 6 | $4.50 \times 10^{-7}$ | 0.229 |
| 5 | $((((12)(34))5)6)$ | 2 | $1.05 \times 10^{-7}$ | 0.054 |
| 6 | $(((12)(34))(56))$ | 8 | $9.06 \times 10^{-7}$ | 0.231 |

orthant taken up by the comb tree topology appears to be the smallest. We predict that these patterns hold for larger numbers of taxa as well.

The Neighbor-Joining (NJ) algorithm (Algorithm 5.4.6) is similar to UPGMA but returns arbitrary tree metrics that may be realized by unrooted trees. See Figure 1.8 for a picture of two unrooted phylogenetic [5]-trees. Like every distance-based phylogenetic reconstruction method, NJ partitions $\mathbb{R}^{n(n-1)/2}$ into a family of regions indexed by the combinatorial type of tree returned by the algorithm. The selection criteria and distance recalculations in the NJ algorithm may all be expressed as linear combinations of the original dissimilarity map inputs, so as is the case with UPGMA, each NJ region is a union of polyhedral cones. NJ does not return rooted trees so there is no notion of a rank function, but as with UPGMA the union is taken over all orderings on the internal vertices that correspond to steps of the NJ algorithm. The inequalities determined by steps in the NJ algorithm determine the $\mathcal{H}$-representation of the NJ cones.

However, a complete description for all $n$ of the extreme rays of the NJ cones is unknown. Such a description would allow one to establish an analogue of Theorem 4.3.7 for NJ. Such a result would be desirable, not only to facilitate more analyses similar to those undertaken in this chapter, but also to aid in the study of other geometric problems associated to NJ such as the problems we will encounter in Chapter 5. One issue with the NJ cones that complicates the situation is that they have a non-trivial lineality space (Definition 1.2.7), so there is more than one representation for an extreme ray. In the paper [32] of Kord Eickmeyer and Ruriko Yoshida, a detailed study of the extreme rays for $n = 5$ is undertaken, in which the extreme rays of a cone are considered as projections onto the orthogonal complement of the lineality space of the cone.

Table 4.4: Spherical volumes of UPGMA regions $\mathcal{P}(T)$ for $n = 7$.

| | Tree | # Chains | Volume | Fraction of Orthant |
|---|---|---|---|---|
| 1 | $((((((12)3)4)5)6)7)$ | 1 | $2.75 \times 10^{-13}$ | 0.0050 |
| 2 | $(((((12)3)4)5)(67))$ | 5 | $4.82 \times 10^{-12}$ | 0.0435 |
| 3 | $(((((12)3)4)(56))7)$ | 4 | $6.32 \times 10^{-12}$ | 0.0570 |
| 4 | $((((12)3)4)((56)7))$ | 10 | $1.95 \times 10^{-11}$ | 0.1762 |
| 5 | $(((((12)3)(45))6)7)$ | 3 | $4.45 \times 10^{-12}$ | 0.0402 |
| 6 | $((((12)3)(45))(67))$ | 15 | $5.72 \times 10^{-11}$ | 0.2581 |
| 7 | $((((12)3)((45)6))7)$ | 6 | $1.66 \times 10^{-11}$ | 0.0747 |
| 8 | $(((12)3)((45)(67)))$ | 20 | $9.00 \times 10^{-11}$ | 0.2030 |
| 9 | $(((((12)(34))5)6)7)$ | 2 | $1.73 \times 10^{-12}$ | 0.0078 |
| 10 | $((((12)(34))5)(67))$ | 10 | $2.63 \times 10^{-11}$ | 0.0593 |
| 11 | $((((12)(34))(56))7)$ | 8 | $3.33 \times 10^{-11}$ | 0.0753 |

# Chapter 5

# Distance-Based Phylogenetic Methods Near a Polytomy

## 5.1 Introduction

Recall that any distance-based phylogenetic reconstruction method $f$ (Definition 1.4.16) takes a dissimilarity map (Definition 1.4.8) as an input and outputs a tree metric (Definition 1.4.9). So every such method $f$ partitions the set of all dissimilarity maps $\mathbb{R}^{n(n-1)/2}$ into regions

$$\{C(T) : f(x) \text{ is a tree metric realized by the combinatorial tree } T \text{ for all } x \in C(T)\}.$$

We can then compare the regions $C(T)$ for different methods to evaluate their relative performance. Among the most intuitively appealing distance-based phylogenetic methods is the *least-squares phylogeny* (LSP):

**Problem 5.1.1.** The least-squares phylogeny problem asks, for a given dissimilarity map $\delta$, what is the tree metric $d$ that minimizes the ordinary Euclidean distance given by the formula

$$\sqrt{\sum_{x,y \in X} (\delta(x,y) - d(x,y))^2}.$$

William Day showed that the least-squares phylogeny problem is NP-hard [28]. Accordingly, many distance-based phylogenetic algorithms have been developed which attempt to build up the tree piece by piece while locally optimizing the Euclidean distance at each step. Two popular agglomerative distance-based methods designed according to this philosophy are the Unweighted Pair-Group Method with Arithmetic Mean, or UPGMA (Algorithm 4.2.1), which we studied in Chapter 4, and the Neighbor-Joining algorithm, or NJ (Algorithm 5.4.6). Both UPGMA and NJ run in polynomial time. As LSP is NP-hard, UPGMA and NJ cannot solve LSP exactly.

So it is natural to ask: how well do these distance-based algorithms perform when attempting to solve the LSP problem? Under what circumstances do distance-based heuristics return the same combinatorial tree as the least-squares phylogeny?

A motivation for the study conducted in this chapter is the following well-known consistency result of Kevin Atteson [4].

**Theorem 5.1.2** (Atteson, 1999). *Let d be a tree metric such that the smallest edge weight in a binary tree realization $(T, w)$ of d is $w(e) > 0$. Suppose a dissimilarity map $\delta \in \mathbb{R}^{n(n-1)/2}$ satisfies $\|\delta - d\|_\infty < w(e)/2$. Then NJ $(\delta) = \hat{d}$, where $\hat{d}$ and d are realized by trees with the same tree topology.*

Theorem 5.1.2 says the following: if $\delta$ is a dissimilarity map which is sufficiently close to some tree metric $d$ realized by a binary tree all of whose branch lengths are bounded away from zero, then NJ applied to $\delta$ returns a tree with the same combinatorial type as $d$. In other words, if our data gives us an input that is close enough in the input space to the correct tree shape, NJ will correctly infer the evolutionary history.

Recall that the Euclidean norm is equivalent to the $\infty$-norm, or in other words for $x \in \mathbb{R}^d$, $\|x\|_\infty \leq \|x\|_2 \leq \sqrt{d}\|x\|_\infty$. Therefore, sufficiently close to an input $\delta$, NJ gives a tree shape consistent with LSP when all edge lengths of a binary tree are bounded away from zero. But, if we allow edge lengths to be equal to zero, this is equivalent to collapsing an edge in a binary tree to obtain a non-binary tree, which results in a *polytomy* vertex (Definition 1.4.3).

In a rooted tree, a polytomy represents a speciation event where many different species were produced. Polytomies routinely arise in phylogenetic inference from collections of species for which there is not enough data to decide which sequence of binary events is most relevant. This leads us to the main question of study in this chapter:

**Problem 5.1.3.** How do the distance based-heuristics UPGMA and NJ compare to LSP when the true tree metric has a polytomy?

The idea of comparing two distanced-based methods using polyhedral geometry already appears in the paper [31] of Kord Eickmeyer, Peter Huggins, Lior Pachter, and Ruriko Yoshida as well as the paper [43] of David Haws, Terrell Hodge, and Yoshida. In both of these papers, Neighbor-Joining is compared to the Balanced Minimum Evolution (BME) criterion. In the case of the distance-based heuristics UPGMA and NJ the resulting regions in the partition of the input space $\mathbb{R}^{n(n-1)/2}$ are polyhedral cones because the decision criteria of both of these algorithms are represented by families of linear inequalities which determine an $\mathcal{H}$-representation of the cones.

For LSP, the regions are potentially more complicated semialgebraic sets, or solutions to families of polynomial inequalities. While we do not yet know a complete description of the

regions induced by LSP, a local analysis of the performance of LSP and distance-based heuristics near a polytomy can be done using polyhedral geometry. The resulting analysis depends heavily on the geometry of phylogenetic tree space (Definition 1.4.15) near tree metrics that contain a polytomy. It is this analysis which comprises the bulk of this chapter. The results in this chapter are joint work with Seth Sullivant, and appear in the paper [26].

This chapter is organized as follows: in Section 5.2, we review basic properties of tree space, including a description of the different cones in the standard decomposition. We provide the description of both tree space $\mathcal{T}_n$ and equidistant tree space $\mathcal{ET}_n$. Section 5.3 contains a detailed analysis of the local geometry of tree space near tree metrics called *tritomies* that are realized by trees which have a polytomy with three associated binary resolutions. In particular, for both equidistant and ordinary tree metrics, the local geometry depends only on the sizes of the daughter clades (Definition 1.4.5) around the tritomy, and not the particular tree structure of those daughter clades.

In Section 5.4, we apply the results from Section 5.3 to understand the local geometry of the decompositions induced by LSP and UPGMA near tree metrics that contain a tritomy. We also explain why these results imply that UPGMA poorly matches LSP in some circumstances, and we discuss computational evidence towards the study of NJ from this perspective. Section 5.5 contains concluding remarks primarily about the possibility of extending results for NJ.

## 5.2 Tree space

Our analysis of the behavior of phylogenetic algorithms near a tree metrics containing a polytomy depends heavily on the geometry of the spaces of tree metrics $\mathcal{T}_n$ and equidistant tree metrics $\mathcal{ET}_n$ on $n$ leaves defined in Definition 1.4.15. Both $\mathcal{T}_n$ and $\mathcal{ET}_n$ are polyhedral fans (Definition 1.2.8). The fan $\mathcal{T}_n$ has one maximal cone for each unrooted binary tree. The space of ultrametrics $\mathcal{ET}_n$ has one maximal cone for each rooted binary tree. The extreme rays of these maximal cones are known in both cases. The space $\mathcal{ET}_3$ is shown in Figure 5.1.

**Definition 5.2.1.** For each $i \neq j \in X$, recall $e_{ij} \in \mathbb{R}^{n(n-1)/2}$ denotes the dissimilarity map such that $e_{ij}(i,j) = 1$ and $e_{ij}(x,y) = 0$ for all other pairs $x, y$. Let $A_1, A_2, \ldots, A_k$ be a collection of disjoint subsets of $X$. Define the dissimilarity map $\delta_{A_1|A_2|\cdots|A_k}$

$$\delta_{A_1|A_2|\cdots|A_k} = \sum_{ij} e_{ij}$$

where the sum ranges over all unordered pairs $(i, j)$ such that $i$ and $j$ belong to different blocks.

In the special case where $A|B$ is a partition of $X$, $A|B$ is usually called a *split*. The resulting dissimilarity map $\delta_{A|B}$ is called a *cut-semimetric* or *split-psuedometric*. Each edge in a tree $T$
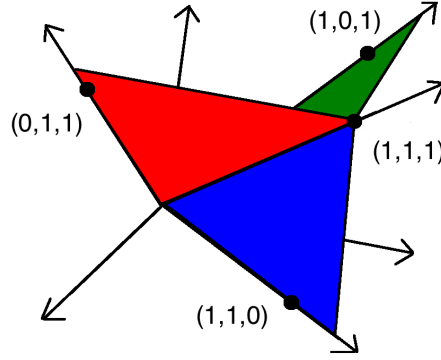
Figure 5.1: The space $\mathcal{ET}_3$ with labeled extreme rays.

induces a split of the leaves of $T$ obtained from the partition of the leaves that arises from removing the indicated edge. The set of all splits implied by a tree $T$ is denoted $\Sigma(T)$. Recall that cone$(Y) \subset \mathbb{R}^d$ is *simplicial* if $Y$ is a linearly independent set of elements of $\mathbb{R}^d$.

**Proposition 5.2.2.** *Let $T$ be a phylogenetic $X$-tree. The set of all tree metrics compatible with $T$ is a simplicial cone, whose extreme rays are the set of vectors $\{\delta_{A|B} : A|B \in \Sigma(T)\}$.*

This is a polyhedral geometry rewording of Theorem 7.1.8 in the book *Phylogenetics* by Charles Semple and Mike Steel [57]. Note that the description from Proposition 5.2.2 holds regardless of whether or not the tree $T$ is binary. In particular, we see that the intersection of all cones associated to a collection of trees corresponds to the cone associated to the tree obtained from a common coarsening of all trees in the given collection.

The cones of the space of equidistant trees $\mathcal{ET}_n$ are not simplicial in general, but they can be subdivided into cones based on ranked trees, which are simplicial. We describe these cones now. Recall that we introduced the lattice $\Pi_n$ of partitions of the set $[n] = \{1, \dots, n\}$ in Example 1.1.4.

**Proposition 5.2.3.** *Let*

$$C = \pi_n \lessdot \pi_{n-1} \lessdot \cdots \lessdot \pi_1$$

*be a maximal chain in $\Pi_n$, corresponding to a ranked phylogenetic tree. The cone of equidistant tree metrics compatible with $C$ is a simplicial cone whose extreme rays are the set of vectors $\{\delta_{\pi_i} : i = 2, \dots, n\}$.*

This is a polyhedral geometry rewording of Theorem 7.2.8 of [57].

**Example 5.2.4.** If an equidistant tree metric $d \in \mathbb{R}^{10}_{\geq 0}$ is compatible with the maximal chain in Example 1.4.7 then $d$ satisfies

$$d(1,2) \leq d(3,4) \leq d(1,3) = d(1,4) = d(2,3) = d(2,4) \leq$$

$$d(1,5) = d(2,5) = d(3,5) = d(4,5)$$

and is in the simplicial cone with extreme rays

$$(1,1,1,1,1,1,1,1,1,1), \quad (0,1,1,1,1,1,1,1,1,1),$$

$$(0,1,1,1,1,1,1,0,1,1), \quad (0,0,0,1,0,0,1,0,1,1)$$

where the coordinates of $\mathbb{R}^{10}$ are labeled with the pairs $i < j \in [5]$ in the lexicographic order.

Note that Proposition 5.2.3 also holds true when working with chains that are not maximal, which correspond to either trees with polytomies or situations where there are ties in the rankings of the internal vertices. These chains correspond to intersections of the maximal cones associated to the maximal chains in the partition lattice $\Pi_n$, which, as we saw in Section 1.4, correspond to rooted, ranked phylogenetic $[n]$-trees.
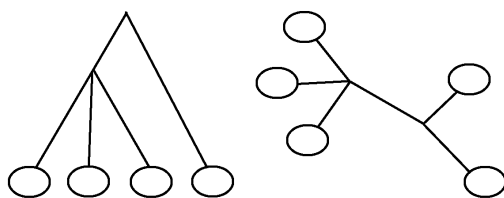


Figure 5.2: Two tritomies, rooted and unrooted.

## 5.3  Geometry of Tree Space Near a Tritomy

The goal of this section is to describe the geometry of tree space near a polytomy, in particular in the special case of tritomies. For rooted trees, a *tritomy* is an internal vertex that has three direct descendants. In an unrooted tree a *tritomy* is an internal vertex with four neighbors. Figure 5.2 shows a rooted tritomy tree and an unrooted tritomy tree.

When we speak of the "geometry of tree space near a tritomy", we mean to describe the geometry of tree space near a generic tree metric that is realized by a tree with a single tritomy and no other polytomies. The set of all such tritomy tree metrics, for a fixed topological structure on the tree $T$, is a polyhedral cone of dimension one less than the dimension of tree space. Let $C_T$ denote this polyhedral cone. The tree $T$ with a single tritomy can be resolved to three binary trees. Denote them $T_1, T_2$, and $T_3$. The polyhedral cone of a tritomy $C_T$ is the intersection of the three *resolution cones* $C_{T_1}$, $C_{T_2}$, and $C_{T_3}$ associated to the three different ways to resolve
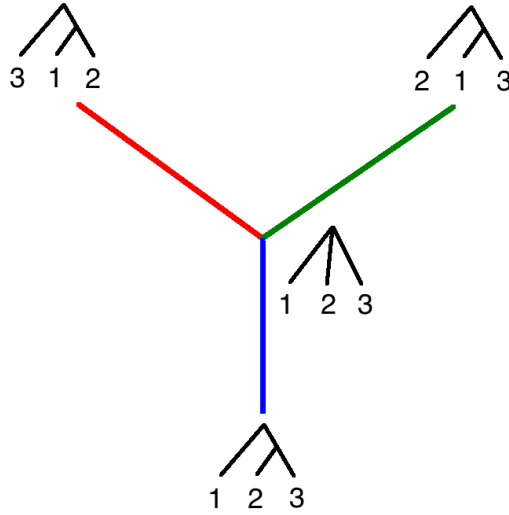
79

Figure 5.3: The fan $K_T$ with labeled cones.

the tritomy tree into a binary tree.

For both equidistant tree space $\mathcal{ET}_n$ and ordinary tree space $\mathcal{T}_n$, the cones $C_T$ and their resolution cones $C_{T_i}$, for $i = 1, 2$, and $3$ satisfy $\dim C_T = \dim C_{T_i} - 1$. This is easily seen by the simplicial structure of the cones $C_T$ for any tree $T$, according to Propositions 5.2.2 and 5.2.3. Hence, locally near a generic point $x$ of $C_T$, tree space looks like $\mathbb{R}^k \times K_T$ where $k = \dim C_T$ and $K_T$ is a one-dimensional polyhedral fan that depends on $T$ but does not depend on $x$. Furthermore, the fan $K_T$ can be chosen to live in a space orthogonal to the span of $C_T$, and span $K_T$ is two-dimensional. Figure 5.3 shows the fan $K_T$ in $(\text{span } C_T)^{\perp}$ when $T$ is the rooted tritomy tree with leaf set $[3]$.

The goal of this section is to describe the structure of that fan $K_T$. The analysis depends on the particular structure of the generators of the various cones involved, and the cases of equidistant tree metrics and arbitrary tree metrics must be handled separately. We treat these cases in Sections 5.3.1 and 5.3.2, respectively.

### 5.3.1    Equidistant Tree Space

In this section we determine the geometry of the fan $K_T$ for a tritomy tree $T$ in equidistant tree space. This tritomy tree has a node with three children. Denote the daughter clades of these children (that is, the set of leaves below each of the children of the tritomy, see Definition 1.4.5) by $A$, $B$, and $C$. Let $T_{AB}$, $T_{AC}$, and $T_{BC}$ denote the three resolution trees, where for example $T_{AB}$ is the binary resolution where $A \cup B$ forms the leaf set of a clade. Note that since all the linear spaces that are involved are the same, instead of working with a fixed tree we can work

with the corresponding rank function and chain in $\Pi_n$, which are derived from the order of agglomeration of subsets of the leaf set induced by the tree topology as illustrated in Example 1.4.7. This is what we will do in this section.

Let

$$K = \pi_n \lessdot \pi_{n-1} \lessdot \cdots \lessdot \pi_{k+1} \lessdot \pi_{k-1} \lessdot \cdots \lessdot \pi_1$$

be the chain corresponding to the polytomy tree. Note that this is a chain in the partition lattice which leaves out an element at the $k$-th level. Here $\pi_{k+1}$ will contain among its blocks $A, B$, and $C$, and $\pi_{k-1}$ will contain the block $A \cup B \cup C$. The resolution trees $T_{AB}$, $T_{AC}$, and $T_{BC}$ correspond to the three ways to add a $\pi_k$ to this sequence which refines $\pi_{k+1}$ and is refined by $\pi_{k-1}$.

We are interested in the linear spaces span $C_K$.

**Lemma 5.3.1.** *For any (not necessarily maximal) chain $K = \pi_r \lessdot \cdots \lessdot \pi_1$, where $\pi_1 = X$, the set of vectors*

$$\{\delta_{\pi_i} - \delta_{\pi_{i-1}}\}_{i=2,\ldots,r}$$

*forms an orthogonal basis for* span $C_K$.

*Proof.* Since $\delta_{\pi_2}, \ldots, \delta_{\pi_r}$ are the extreme rays of the simplicial cone $C_K$, they are linearly independent and hence span the space span $C_K$. We can easily solve for the vectors $\delta_{\pi_2}, \ldots, \delta_{\pi_r}$ given $\delta_{\pi_i} - \delta_{\pi_{i-1}}, i = 2, \ldots, r$ hence span $C_K = $ span $\{\delta_{\pi_i} - \delta_{\pi_{i-1}}\}_{i=2,\ldots,r}$. For all $i \in [r]$, the positions of the ones in $\delta_{\pi_{i-1}}$ are a subset of the positions of the ones in $\delta_{\pi_i}$. This guarantees that $\delta_{\pi_i} - \delta_{\pi_{i-1}}$ and $\delta_{\pi_j} - \delta_{\pi_{j-1}}$ do not have any nonzero entries in the same positions when $i \neq j$. Hence these vectors are orthogonal. Note that $\delta_{\pi_1}$ is the zero vector if we assume that $\pi_1 = X$. $\qquad\square$

The particular structure of the vectors $\delta_{\pi_{i+1}} - \delta_{\pi_i}$ will be useful in what follows.

**Lemma 5.3.2.** *Let $\pi$ and $\tau$ be two set partitions such that $\pi$ is a refinement of $\tau$. Then*

$$\delta_\pi - \delta_\tau = \sum e_{i,j}$$

*where $i, j$ are in different parts of $\pi$ and the same part of $\tau$.*

*Proof.* Trivial from the definition of $\delta_\pi$. $\qquad\square$

Now for each of the resolution cones, for example $C_{T_{AB}}$, there is a unique ray $p_{AB}$ in the fan $K_T$ that is orthogonal to span $C_T$. Since $\dim C_{T_{AB}} = \dim C_T + 1$ and the cones $C_{T_{AB}}$ and $C_T$ are pointed (see the comment after Definition 1.2.7), the ray $p_{AB}$ is unique. The ray $p_{AB}$ is orthogonal to span $C_T$ because we choose the fan $K_T$ to be orthogonal to span $C_T$. We explain how to construct that ray now.

**Lemma 5.3.3.** *Let* $a = |A|$, $b = |B|$, *and* $c = |C|$. *The vector* $p_{AB}$ *is given by*

$$p_{AB} = -\frac{ac + bc}{ab + ac + bc}\delta_{A|B} + \frac{ab}{ab + ac + bc}(\delta_{A|C} + \delta_{B|C}).$$

*Proof.* It suffices to start with any vector $r_{AB} \in \operatorname{span} C_{T_{AB}} \setminus \operatorname{span} C_T$ and project it onto the orthogonal complement of $\operatorname{span} C_T$. We assume the tree $T$ is represented by the chain

$$K = \pi_n \lessdot \pi_{n-1} \lessdot \cdots \lessdot \pi_{k+1} \lessdot \pi_{k-1} \lessdot \cdots \lessdot \pi_1$$

and the tree $T_{AB}$ by the chain

$$K_{AB} = \pi_n \lessdot \pi_{n-1} \lessdot \cdots \lessdot \pi_{k+1} \lessdot \pi_k \lessdot \pi_{k-1} \lessdot \cdots \lessdot \pi_1.$$

For our vector we choose $r_{AB} = \delta_{\pi_k} - \delta_{\pi_{k-1}}$. This vector is clearly not in $\operatorname{span} C_K$ since it involves $\delta_{\pi_k}$. Furthermore, $r_{AB}$ is already orthogonal to all the vectors in the orthogonal basis for $\operatorname{span} C_T$, except for the vector $\delta_{\pi_{k+1}} - \delta_{\pi_{k-1}}$. Hence, we can project the vector $r_{AB}$ onto the complement of the space spanned by $\delta_{\pi_{k+1}} - \delta_{\pi_{k-1}}$; this will be the same as projecting onto the complement of $\operatorname{span} C_T$.

Note that by Lemma 5.3.2 $r_{AB} = \delta_{A \cup B|C} = \delta_{A|C} + \delta_{B|C}$. Similarly $\delta_{\pi_{k+1}} - \delta_{\pi_{k-1}} = \delta_{A|B|C} = \delta_{A|B} + \delta_{A|C} + \delta_{B|C}$. So we want to project $\delta_{A|C} + \delta_{B|C}$ onto the orthogonal complement of $\delta_{A|B} + \delta_{A|C} + \delta_{B|C}$. To find $p_{AB}$ it is enough to compute the component of $\delta_{A|C} + \delta_{B|C}$ that is perpendicular to $\delta_{A|B} + \delta_{B|C} + \delta_{A|C}$, otherwise known as the vector rejection of $\delta_{B|C} + \delta_{A|C}$ from $\delta_{A|B} + \delta_{B|C} + \delta_{A|C}$. Due to the orthogonality of the vectors $\delta_{A|B}, \delta_{A|C}$, and $\delta_{B|C}$ and the fact that, for example, $(\|\delta_{A|B}\|_2)^2 = |A||B|$, we have

$$\frac{(\delta_{A|C} + \delta_{B|C}) \cdot (\delta_{A|B} + \delta_{B|C} + \delta_{A|C})}{(\delta_{A|B} + \delta_{B|C} + \delta_{A|C}) \cdot (\delta_{A|B} + \delta_{B|C} + \delta_{A|C})} =$$

$$\frac{(\|\delta_{B|C}\|_2)^2 + (\|\delta_{A|C}\|_2)^2}{(\|\delta_{A|B}\|_2)^2 + (\|\delta_{B|C}\|_2)^2 + (\|\delta_{A|C}\|_2)^2}$$

So the vector rejection of $\delta_{B|C} + \delta_{A|C}$ from $\delta_{A|B} + \delta_{B|C} + \delta_{A|C}$ becomes

$$\delta_{A|C} + \delta_{B|C} - \frac{ac + bc}{ab + bc + ac}(\delta_{A|B} + \delta_{B|C} + \delta_{A|C}) =$$

$$-\frac{ac + bc}{ab + ac + bc}\delta_{A|B} + \frac{ab}{ab + ac + bc}(\delta_{A|C} + \delta_{B|C}).$$

$\square$

This explicit formula for $p_{AB}$ implies that $\operatorname{span} K_T$ is two-dimensional:

**Corollary 5.3.4.** *The space* span $K_T$ *is two-dimensional.*

*Proof.* The formulae for $p_{AB}$, $p_{AC}$, and $p_{BC}$ given by Lemma 5.3.3 show that $p_{AB}$ and $p_{AC}$ are not parallel, so that span $K_T$ is at least two-dimensional, but

$$p_{AB} + p_{AC} + p_{BC} = \mathbf{0},$$

where $\mathbf{0}$ denotes the zero vector in $\mathbb{R}^{n(n-1)/2}$. So span $K_T$ is exactly two-dimensional. $\square$

Corollary 5.3.4 also follows from the fact that the set of all equidistant tree metrics is a tropical variety in $\mathbb{R}^{n(n-1)/2}$, as shown by David Speyer and Bernd Sturmfels in [62].

**Theorem 5.3.5.** *The angle between the cones $C_{T_{AC}}$ and $C_{T_{BC}}$ is*

$$\arccos\left(\frac{-c}{\sqrt{(a+c)(b+c)}}\right).$$

*Proof.* We must calculate the angle between the vectors $p_{AC}$ and $p_{BC}$. This is

$$\arccos\left(\frac{p_{AC} \cdot p_{BC}}{\|p_{AC}\|_2 \|p_{BC}\|_2}\right).$$

Now

$$p_{AC} = -\frac{ab+bc}{ab+ac+bc}\delta_{A|C} + \frac{ac}{ab+ac+bc}(\delta_{A|B} + \delta_{B|C})$$

and

$$p_{BC} = -\frac{ab+ac}{ab+ac+bc}\delta_{B|C} + \frac{bc}{ab+ac+bc}(\delta_{A|B} + \delta_{A|C}).$$

Thus, $p_{AC} \cdot p_{BC}$ is given by

$$-\frac{ab+bc}{ab+ac+bc} \cdot \frac{bc}{ab+ac+bc} \cdot \|\delta_{A|C}\|_2^2$$

$$+\frac{ac}{ab+ac+bc} \cdot \frac{bc}{ab+ac+bc} \cdot \|\delta_{A|B}\|_2^2$$

$$-\frac{ac}{ab+ac+bc} \cdot \frac{ab+ac}{ab+ac+bc} \cdot \|\delta_{B|C}\|_2^2$$

$$= \frac{-a^2bc}{ab+bc+ac}.$$

Similar calculations show that

$$\|p_{AC}\|_2 = \sqrt{\frac{abc(a+c)}{ab+ac+bc}} \quad\text{and}\quad \|p_{BC}\|_2 = \sqrt{\frac{abc(b+c)}{ab+ac+bc}}.$$

Combining these pieces produces the formula in the Theorem. $\square$

### 5.3.2 Tree Space

In this section we determine the geometry of the fan $K_T$ for a tritomy tree $T$ in unrooted tree space $\mathcal{T}_n$. The approach is similar to the analysis for equidistant tree space, but the structure of tree space is more complicated. In particular, finding an orthogonal basis for the space spanned by the rays of the intersection cone for a tritomy is less straightforward.

Recall that a tritomy $p$ in an unrooted tree is an internal vertex of degree four. The edges adjacent to $p$ induce a four-way set partition $A|B|C|D$ of $[n]$. Let $T_{AB}$ denote the resolution tree in which there is an edge inducing the split $A \cup B | C \cup D$. Note that $T_{AB} = T_{CD}$. So there are three resolutions $T_{AB}, T_{AC}$, and $T_{AD}$ of $T$. For the remainder of this section, let $a = |A|$, $b = |B|$, $c = |C|$, and $d = |D|$. Let $r_{AB} = \delta_{A \cup B | C \cup D}$. Note that $r_{AB}$ and the projection $p_{AB}$ of $r_{AB}$ onto $(\operatorname{span} C_T)^\perp$ are different objects than the $r_{AB}$ and $p_{AB}$ that we defined in the case of equidistant tree space, but we use the same notation to illustrate the parallel roles that $r_{AB}$ and $p_{AB}$ play in this discussion.

**Lemma 5.3.6.** *Let $T$ be an unrooted tree with a tritomy and corresponding partition $A|B|C|D$ of $[n]$. Then $r_{AB} \in \operatorname{span} C_{T_{AB}} \setminus \operatorname{span} C_T$, and $\dim \operatorname{span} C_T = 2n - 4$.*

*Proof.* By Proposition 5.2.2 each extreme ray of $C_{T_{AB}}$ (equivalently, $T$) comes from a split induced by an edge of $T_{AB}$ (equivalently, an edge of $T$). A binary unrooted tree on $n$ leaves has $2n - 3$ edges. So the cone $C_{T_{AB}}$ has $2n - 3$ rays, one for each internal edge of the tree $T_{AB}$. By contracting the edge in $T_{AB}$ that induces the split $A \cup B | C \cup D$ we obtain $T$. Therefore $C_T$ has $2n - 4$ extreme rays that correspond to the $2n - 4$ internal edges of $T$. Since $C_T$ is simplicial, $\dim \operatorname{span} C_T = 2n - 4$. $\square$

The projections $p_{AB}$, $p_{AC}$, and $p_{AD}$ of $r_{AB}$, $r_{AC}$, and $r_{AD}$ onto $(\operatorname{span} C_T)^\perp$ are the maximal cones in the fan $K_T$. As in the previous section, we use an orthogonal basis of $\operatorname{span} C_T$ to simplify the necessary calculations.

Let $A'$ denote the complement of $A$ in the set $[n]$. The vectors in the set

$$\mathcal{U} = \{\delta_{A|A'}, \delta_{B|B'}, \delta_{C|C'}, \delta_{D|D'}\}$$

are extreme rays of $T$. The elements of $\mathcal{U}$ correspond to the four edges in $T$ adjacent to the tritomy $p$. We show in the next Lemma that to calculate $p_{AB}$ it is sufficient to calculate the projection of $r_{AB}$ onto $(\operatorname{span} \mathcal{U})^\perp$. First we require some additional notation: let $e = (u, v)$ be an edge of $T$ not adjacent to $p$ where $v$ is the internal vertex of $T$ on the path to $p$ from $e$. Let $e' = (w, v)$ be the unique edge in $T$ satisfying the conditions $(i)$ $e \neq e'$ and $(ii)$ $w$ appears on the path from $v$ to $p$ in $T$ (note that it is possible that $w = p$). Let $A_e|B_e$ be the split of $[n]$ induced by $e$ and let $A_{e'}|B_{e'}$ be the split of $[n]$ induced by $e'$. Note $A_e \subsetneq A_{e'}$. Let $a_e = |A_e|$ and

$a_{e'} = |A_{e'}|$. Let $\mathcal{V}$ be the set

$$\left\{ \delta_{A_e|B_e} - \frac{a_e}{a_{e'}} \delta_{A_{e'}|B_{e'}} : p \notin e = (u,v), e' \text{ satisfies } (i), (ii) \right\}.$$

**Lemma 5.3.7.** *Every vector in $\mathcal{V}$ is orthogonal to $r_{AB}, r_{AC}$, and $r_{AD}$ and $\mathcal{U} \cup \mathcal{V}$ is a basis for* span $C_T$.

*Proof.* Since $T$ has exactly one tritomy, $T$ has $2n-4$ edges. When $n = 4$, $2n - 4 = 4$. In this case $|\mathcal{U}| = \dim \text{span } C_T$, and $\mathcal{U}$ is a basis for span $C_T$. So, assume $n > 4$: then $\mathcal{V}$ is not empty because we can find edges $e$ and $e'$ satisfying $p \notin (u,v) = e$ and $e'$ satisfying conditions $(i)$ and $(ii)$. We will first show that each element of $\mathcal{V}$ is orthogonal to $r_{AB}, r_{AC}$, and $r_{AD}$. Let $\nu \in \mathcal{V}$, then

$$\nu = \delta_{A_e|B_e} - \frac{a_e}{a_{e'}} \delta_{A_{e'}|B_{e'}}.$$

Note that $A_{e'}$ is contained in one of $A, B, C$, and $D$. Without loss of generality we may assume that $A_{e'} \subset A$. Then it follows directly from the structure of the summands in the vector $\nu$ that

$$r_{AB} \cdot \nu = (a_{e'})(c+d)\left(-\frac{a_e}{a_{e'}}\right) + a_e(c+d) = 0.$$

Similar calculations show that $r_{AC}$ and $r_{AD}$ are also orthogonal to $\nu$.

We obtain $2n - 8$ vectors in $\mathcal{V}$ because there are $(2n-4) - 4$ edges in $T$ that do not induce vectors in $\mathcal{U}$. So $|\mathcal{U} \cup \mathcal{V}| = 2n-4$. Since $\mathcal{U} \cup \mathcal{V}$ is comprised of vectors that are linear combinations of split-pseudometrics, span $\mathcal{U} \cup \mathcal{V} \subset$ span $C_T$. The set $\mathcal{U} \cup \mathcal{V}$ is also linearly independent since it can be seen as an upper triangular transformation of the set of extreme rays of $C_T$, which were independent. Thus $\mathcal{U} \cup \mathcal{V}$ is a basis for span $C_T$. $\square$

Due to the structure of the vectors $r_{AB}$ and the elements of $\mathcal{U}$, $p_{AB}$ is constant on the coordinates for each $\delta_{U|V}$. So we can write

$$p_{AB} = \sum_{\{U,V\} \in \binom{\{A,B,C,D\}}{2}} w(AB)_{U|V} \cdot \delta_{U|V}. \tag{5.1}$$

The coefficients $w(AB)_{U|V}$ will facilitate computation of dot products and 2-norms.

**Theorem 5.3.8.** *The angle between the cones $C_{T_{AB}}$ and $C_{T_{AC}}$ is*

$$\arccos\left(-\frac{bc + ad}{\sqrt{(a+b)(a+c)(b+d)(c+d)}}\right).$$

*Proof.* By Lemma 5.3.7, to find $p_{AB}, p_{AC}$, and $p_{AD}$ it is sufficient to calculate the projection of $r_{AB}, r_{AC}$, and $r_{AD}$ onto $(\text{span } \mathcal{U})^\perp$. We will find $p_{AB}$ by calculating the coefficients

85

$w(AB)_{U|V}$. First, we construct a matrix $M_{AB}$ shown in (5.2) as follows: the rows of $M_{AB}$ are indexed by the vectors $(\delta_{A|A'}, \delta_{B|B'}, \delta_{C|C'}, \delta_{D|D'})$, and the columns are indexed by the vectors $(r_{AB}, \delta_{A|A'}, \delta_{B|B'}, \delta_{C|C'}, \delta_{D|D'})$. $(M_{AB})_{i,j}$ is obtained (up to row operations) by taking the dot product of the vector indexing row $i$ and column $j$.

$$
\begin{bmatrix}
c+d & b+c+d & b & c & d \\
c+d & a & a+c+d & c & d \\
a+b & a & b & a+b+d & d \\
a+b & a & b & c & a+b+c
\end{bmatrix}
\tag{5.2}
$$

Next, let $K_{AB}$ be the matrix in (5.3) with rows and columns indexed as shown.

| | $r_{AB}$ | $\delta_{A|A'}$ | $\delta_{B|B'}$ | $\delta_{C|C'}$ | $\delta_{D|D'}$ |
|---|---|---|---|---|---|
| $\delta_{A|B}$ | 0 | 1 | 1 | 0 | 0 |
| $\delta_{A|C}$ | 1 | 1 | 0 | 1 | 0 |
| $\delta_{A|D}$ | 1 | 1 | 0 | 0 | 1 |
| $\delta_{B|C}$ | 1 | 0 | 1 | 1 | 0 |
| $\delta_{B|D}$ | 1 | 0 | 1 | 0 | 1 |
| $\delta_{C|D}$ | 0 | 0 | 0 | 1 | 1 |

$$\tag{5.3}$$

Let $\sigma_{AB} = \sigma_1 r_{AB} + \sigma_2 \delta_{A|A'} + \sigma_3 \delta_{B|B'} + \sigma_4 \delta_{C|C'} + \sigma_5 \delta_{D|D'}$ be a vector in the null space of the matrix $M_{AB}$. Up to a scalar multiple,

$$
K_{AB}(\sigma_{AB}) =
\begin{bmatrix}
w(AB)_{A|B} \\
w(AB)_{A|C} \\
w(AB)_{A|D} \\
w(AB)_{B|C} \\
w(AB)_{B|D} \\
w(AB)_{C|D}
\end{bmatrix}
=
\begin{bmatrix}
(a+b)cd(c+d) \\
-bd(bc+ad) \\
-bc(ac+bd) \\
-ad(ac+bd) \\
-ac(bc+ad) \\
ab(a+b)(c+d)
\end{bmatrix}.
$$

Then

$$
p_{AB} \cdot p_{AC} = \sum_{\{U,V\} \in \binom{\{A,B,C,D\}}{2}} |U| \cdot |V| \cdot w(AB)_{U|V} \cdot w(AC)_{U|V}
\tag{5.4}
$$

and

$$
\|p_{AB}\|_2 = \sqrt{\sum_{\{U,V\} \in \binom{\{A,B,C,D\}}{2}} |U| \cdot |V| \cdot [w(AB)_{U|V}]^2}.
\tag{5.5}
$$

We use (5.4) and (5.5) to obtain the formulae for the angle measures between the resolution cones. $\qquad \square$

**Corollary 5.3.9.** *The space* span $K_T$ *is two-dimensional when $T$ is unrooted.*

*Proof.* As in the case of rooted trees, we can use the formula in (5.1) to show that the set $\{p_{AB}, p_{AC}\}$ is linearly independent, but the set $\{p_{AB}, p_{AC}, p_{BC}\}$ is linearly dependent. □

As was the case with Corollary 5.3.4, Corollary 5.3.9 also follows from the fact that the set of all arbitrary tree metrics is a tropical variety in $\mathbb{R}^{n(n-1)/2}$, as shown in [62].

## 5.4 Distance-Based Methods Near a Tritomy

In this section, we analyze the performance of distance-based methods around a tritomy using the results about the geometry of tree space from Section 5.3. While in many cases we do not have a complete understanding of the geometry of these decompositions across all of $\mathbb{R}_{\geq 0}^{n(n-1)/2}$, we can describe the geometry in a small neighborhood of a tree metric with a single tritomy. It is this geometry which we explore in the present section.

### 5.4.1 Least-Squares Phylogeny

Let $C_1, \ldots, C_r$ be subsets of $\mathbb{R}^{n(n-1)/2}$. The *Voronoi cell $V_k$* associated with the subset $C_k$ is the set of all points

$$\{x \in \mathbb{R}^{n(n-1)/2} : d(x, C_k) \leq d(x, C_j) \text{ for all } j \neq k\}$$

where $d(x, C_k) = \inf\{||x - a||_2 : a \in C_k\}$. The *Voronoi decomposition* is the subdivision of $\mathbb{R}^{n(n-1)/2}$ into Voronoi cells of the set $\{C_k\}$. When $\{C_k\}$ is the collection of cones associated to all possible combinatorial trees with leaf set $[n]$, the Voronoi cells comprise the subdivision of space induced by the least-squares phylogeny problem.

While the Voronoi decomposition of a finite set of points is well-known to be a polyhedral subdivision of space, the Voronoi decomposition induced by a collection of higher-dimensional polyhedra can be a complicated semi-algebraic decomposition. Hence, the Voronoi decomposition induced by the tree cones is probably not polyhedral. We saw in Section 3 that in a neighborhood of a tree metric $T$ with a single tritomy, tree space has the form $\mathbb{R}^k \times K_T$, where $k = \dim \operatorname{span} C_T$ and $K_T$ is a one-dimensional fan with three rays that sits naturally inside a two-dimensional linear space span $K_T$. In this setting it is easy to describe the Voronoi decomposition.

**Proposition 5.4.1.** *Let $T$ be a tree metric with a single tritomy in $\mathbb{R}_{\geq 0}^{n(n-1)/2}$ with local tree space $\mathbb{R}^k \times K_T$. The boundary between the Voronoi cells for the resolution cones $C_{T_{AB}}$ and $C_{T_{AC}}$ is completely determined by the angle bisector in* span $K_T$ *between $p_{AB}$ and $p_{AC}$.*

*Proof.* The Euclidean distance between the cones $C_{T_{AB}}$ and $C_{T_{AC}}$ is the sum of the distances between them in the two orthogonal spaces span $C_T$ and span $K_T$. Of these two distances, only the distance in the two-dimensional space span $K_T$ is nonzero; this distance is determined by the angle between the maximal cones in the one-dimensional polyhedral fan $K_T$. The set of all points in the plane span $K_T$ equidistant between two vectors emanating from the origin is the bisector of the angle between the two vectors. $\qquad\square$

Proposition 5.4.1 allows us to easily compute the relative size of the Voronoi regions around a polytomy for either equidistant or ordinary tree metrics. The next theorem gives a formula for the boundary between the Voronoi regions.

**Theorem 5.4.2.** *Let $T$ be a ranked, rooted tree with a single tritomy. The boundary in* span $K_T$ *between the Voronoi cells for the resolution cones $C_{T_{AB}}$ and $C_{T_{AC}}$ is spanned by the vector*

$$\frac{p_{AB}}{\sqrt{a+b}} + \frac{p_{AC}}{\sqrt{a+c}}.$$

*Proof.* By Proposition 5.4.1 the boundary we wish to compute is given by the angle bisector in span $K_T$ between $p_{AB}$ and $p_{AC}$, which is spanned by the normalized average of the two vectors. By Lemma 5.3.3 we have

$$\|p_{AB}\|_2 = \sqrt{\frac{abc(a+b)}{ab+ac+bc}} \text{ and } \|p_{AC}\|_2 = \sqrt{\frac{abc(a+c)}{ab+ac+bc}}$$

Therefore

$$\frac{p_{AB}}{\|p_{AB}\|_2} + \frac{p_{AC}}{\|p_{AC}\|_2} = \sqrt{\frac{ab+ac+bc}{abc}} \left( \frac{p_{AB}}{\sqrt{a+b}} + \frac{p_{AC}}{\sqrt{a+c}} \right),$$

which is a multiple of the vector that we claimed spans the boundary of the Voronoi cells for the two cones in span $K_T$. $\qquad\square$

**Theorem 5.4.3.** *Let $T$ be an unrooted tree with a single tritomy. The boundary of the Voronoi cell in* span $K_T$ *between the resolution cones $C_{T_{AB}}$ and $C_{T_{AC}}$ is spanned by the vector*

$$\frac{p_{AB}}{\sqrt{(a+b)(c+d)}} + \frac{p_{AC}}{\sqrt{(a+c)(b+d)}}.$$

*Proof.* We use the fact that $p_{AB}$ and $p_{AC}$ are constant on the vectors $\delta_{U|V}$ for $\{U,V\} \in \binom{\{A,B,C,D\}}{2}$ and the formulae for the coeffcients $w(A,B)_{U|V}$ derived in the proof of Theorem 5.3.8 to calculate the 2-norms of $p_{AB}$ and $p_{AC}$. Up to an identical polynomial $f$ in the variables $a, b, c$ and $d$, we have

$$\|p_{AB}\|_2 = f \cdot \sqrt{(a+b)(c+d)} \qquad\qquad (5.6)$$

and
$$\|p_{AC}\|_2 = f \cdot \sqrt{(a+c)(b+d)}. \tag{5.7}$$

As in Theorem 5.4.2 we know that the angle bisector between $p_{AB}$ and $p_{AC}$ gives the boundary of the Voronoi cell in span $K_T$. By (5.6) and (5.7) the angle bisector is a multiple of

$$\frac{p_{AB}}{\sqrt{(a+b)(c+d)}} + \frac{p_{AC}}{\sqrt{(a+c)(b+d)}}.$$

$\square$

### 5.4.2  UPGMA Regions Near a Polytomy

In this section we show that in some circumstances, UPGMA (Algorithm 4.2.1) fails to correctly identify the least-squares phylogeny. The occurrence and severity of this failure depends entirely on the relative sizes of the daughter clades $A$, $B$, and $C$ of the tritomy.

Recall that Equation 4.1 tells us that if the blocks $A, B \subset [n]$ are joined in step $i$ of UPGMA the distance recalculation implies

$$\delta^i(A, B) = \frac{1}{|A||B|} \sum_{x \in A, y \in B} \delta^n(x, y) =$$

$$\frac{1}{ab} \sum_{x \in A, y \in B} \delta(x, y),$$

a formula which is useful in the proof of the next proposition:

**Proposition 5.4.4.** *Let $T$ be a ranked, rooted tree with a single tritomy. The boundaries between the UPGMA regions in $\mathbb{R}_{\geq 0}^{n(n-1)/2}$ for the resolution cones $C_{T_{AB}}$, $C_{T_{AC}}$, and $C_{T_{BC}}$ are orthogonal to the plane* span $K_T$.

*Proof.* The boundary between the UPGMA regions for the cones $C_{T_{AC}}$ and $C_{T_{BC}}$ is given by the condition

$$\delta^k(A, C) = \delta^k(B, C)$$

which translates into the following linear condition on the original dissimilarity map

$$\frac{1}{ac} \sum_{i \in A, j \in C} \delta(i, j) = \frac{1}{bc} \sum_{i \in B, j \in C} \delta(i, j).$$

This hyperplane has normal vector

$$\frac{1}{ac} \delta_{A|C} - \frac{1}{bc} \delta_{B|C}.$$

Now

$$-\frac{1}{ac}p_{AC} = \frac{1}{ab + ac + bc}\left(-\delta_{A|B} - \delta_{B|C} + \frac{ab + bc}{ac}\delta_{A|C}\right)$$

and

$$\frac{1}{bc}p_{BC} = \frac{1}{ab + ac + bc}\left(\delta_{A|B} + \delta_{A|C} - \frac{ab + ac}{bc}\delta_{B|C}\right).$$

So

$$-\frac{1}{ac}p_{AC} + \frac{1}{bc}p_{BC} =$$

$$\frac{1}{ab + ac + bc}\left(\frac{ab + bc}{ac} + 1\right)\delta_{A|C}$$

$$-\frac{1}{ab + ac + bc}\left(\frac{ab + ac}{bc} + 1\right)\delta_{B|C}$$

$$= \frac{1}{ac}\delta_{A|C} - \frac{1}{bc}\delta_{B|C}.$$

Thus, the normal vector for the boundary between UPGMA regions for the cones $C_{T_{AC}}$ and $C_{T_{BC}}$ is in span $K_T$, and the boundary is orthogonal to span $K_T$. The calculation is the same for the other two pairs of cones. □

**Theorem 5.4.5.** *The boundary between the UPGMA cells for the resolution tree topologies $T_{AC}$ and $T_{BC}$ in span $K_T$ is $-p_{AB}$.*

*Proof.* Since span $K_T$ is two-dimensional and the boundaries between the UPGMA regions for the resolutions $T_{AB}$, $T_{AC}$, and $T_{BC}$ are orthogonal to $K_T$ by Proposition 5.4.4, it suffices to find a vector $\omega \in$ span $K_T$ that satisfies

$$\frac{1}{ac}\sum_{i \in A, j \in C} \omega_{i,j} = \frac{1}{bc}\sum_{k \in B, \ell \in C} \omega_{k,\ell} \leq \frac{1}{ab}\sum_{m \in A, n \in B} \omega_{m,n}.$$

Any such vector will span the boundary of the UPGMA cells. We have

$$\frac{1}{ac}\sum_{i \in A, j \in C} -(p_{AB})_{i,j} =$$

$$\frac{1}{bc}\sum_{k \in B, \ell \in C} -(p_{AB})_{k,\ell} = -\frac{ab}{ab + ac + bc}$$

while

$$\frac{1}{ab}\sum_{m \in A, n \in B} -(p_{AB})_{m,n} = \frac{ac + bc}{ab + ac + bc}$$

So $-p_{AB}$ satisfies the required condition. □

### 5.4.3   UPGMA and LSP Cells

In this section we discuss how results from Sections 5.3 and 5.4 show that UPGMA poorly matches LSP in some circumstances. The geometry of the fan $K_T$ and the UPGMA cells in span $K_T$ for equidistant trees depends entirely on the size of the daughter clades (see Definition 1.4.5) $A, B,$ and $C$ of the tritomy. Consequentially, the quality of the performance of UPGMA near a tree metric with a tritomy depends on how similar in size the daughter clades are. When $a = b = c$, the UPGMA and LSP regions near a tritomy are the same. But if one of the daughter clades becomes much larger or much smaller than the other two daughter clades, UPGMA does a poorer job of identifying the LSP. We also use results from Section 5.4 to show that NJ poorly matches LSP in specific examples for small numbers of taxa.

We can use our theorems about the geometry of span $K_T$ to investigate the relative size of the UPGMA and Voronoi cells as $a, b,$ and $c$ vary. By Theorem 5.3.5 the angle between $C_{T_{AC}}$ and $C_{T_{BC}}$ is

$$\arccos\left(\frac{-c}{\sqrt{(a+c)(b+c)}}\right)$$

and this is also the angle measure of the UPGMA region associated with the cone $C_{T_{AB}}$. By the angle bisector argument, we see that the angle measure of the LSP region associated to the tree $T_{AB}$ near the tritomy will be:

$$\frac{1}{2}\arccos\left(\frac{-a}{\sqrt{(a+b)(a+c)}}\right) +$$

$$\frac{1}{2}\arccos\left(\frac{-b}{\sqrt{(a+b)(b+c)}}\right).$$

When $c \gg a \approx b$, the angle for the UPGMA region approaches $\pi$ whereas the angle for the LSP region approaches $\pi/2$. Conversely, when $a \approx b \gg c$ the angle for the UPGMA region approaches $\pi/2$ whereas the angle for the LSP region approaches $3\pi/4$. Tables 5.1 and 5.2 compare the sizes of the various regions for differing values of $a, b,$ and $c$. We display the sizes as the percentage of the total amount of the local volume around the tritomy that corresponds to the UPGMA or LSP region for the cone $C_{T_{AB}}$. These tables show that while the convergence to the limiting values is slow, already for small values of $a, b,$ and $c$ there may be significant discrepancies between the region sizes for UPGMA and LSP.

Figures 5.4 and 5.5 illustrate the geometry of this phenomenon for the two extreme cases $c \gg a \approx b$, and $a \approx b \gg c$. In both figures the fan $K_T$ is black, the vector $p_{AB}$ is labeled with the pair $AB$, LSP boundaries are dark gray, and UPGMA boundaries are light gray. Note that when $c \gg a \approx b$, UPGMA overestimates the size of the LSP region for $C_{T_{AB}}$. When $a \approx b \gg c$,

Table 5.1: Region sizes for $C_{T_{AB}}$ when $c \gg a = b$.

| a | b | c | UPGMA | LSP |
|---|---|---|---|---|
| 1 | 1 | 1 | 33.3333 | 33.3333 |
| 1 | 1 | 2 | 36.6139 | 31.693 |
| 1 | 1 | 4 | 39.7583 | 30.1209 |
| 1 | 1 | 8 | 42.4261 | 28.787 |
| 1 | 1 | 16 | 44.5139 | 27.7431 |
| 1 | 1 | $2^{10}$ | 49.297 | 25.3515 |
| 1 | 1 | $2^{20}$ | 49.978 | 25.011 |

Table 5.2: Region sizes for $C_{T_{AB}}$ when $a = b \gg c$.

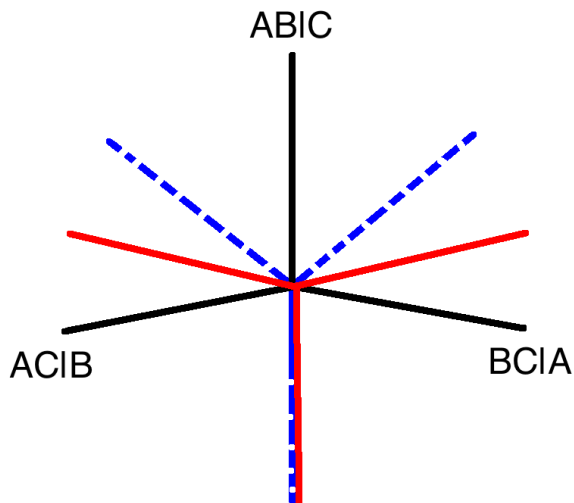| a | b | c | UPGMA | LSP |
|---|---|---|---|---|
| 1 | 1 | 1 | 33.3333 | 33.3333 |
| 2 | 2 | 1 | 30.4086 | 34.7957 |
| 4 | 4 | 1 | 28.2046 | 35.8977 |
| 8 | 8 | 1 | 26.7721 | 36.6139 |
| 16 | 16 | 1 | 25.9367 | 37.0317 |
| $2^{10}$ | $2^{10}$ | 1 | 25.0155 | 37.4923 |
| $2^{20}$ | $2^{20}$ | 1 | 25.0001 | 37.4999 |

Figure 5.4: The case $c \gg a \approx b$. The cones in the fan $K_T$ are labeled with the tritomy resolutions, LSP boundaries are dashed, and UPGMA boundaries are solid.

UPGMA underestimates the size of the LSP region for $C_{T_{AB}}$.

### 5.4.4 LSP Cells and Local NJ Behavior

The Neighbor-Joining (NJ) algorithm, due to Naruya Saitou and Masatoshi Nei [55] is a distance-based reconstruction method that returns an unrooted tree $T$ and a tree metric realized by $T$. Both the selection criterion (known as the "$Q$-criterion") and distance recalculation are linear combinations of the original input coordinates. Therefore, as in the case of UPGMA, NJ divides the input space $\mathbb{R}_{\geq 0}^{n(n-1)/2}$ into a family of polyhedral cones studied in [31] and [43]. In practice, these cones have a simpler description when the non-negativity constraint is relaxed, but these considerations are beyond the scope of the analysis in this chapter, and as with UPGMA, only inputs with positive entries are commonly used by biologists. So, for convenience we will continue to discuss inputs that lie only in the region $\mathbb{R}_{\geq 0}^{n(n-1)/2}$ of $\mathbb{R}^{n(n-1)/2}$.

As discussed in Section 4.6, a complete combinatorial description of the NJ cones remains unknown, and we do not have a closed description of the local geometry of the NJ regions around a tritomy. However, by running NJ on points sampled uniformly from the surface of a small sphere around a tritomy, we can obtain an empirical estimate of the local relative size of NJ regions for small numbers of taxa.

For unrooted tree metrics, the case of interest is when $a$ and $b$ are larger than $c$ and $d$: if $a = b = c$ and $d$ is larger or smaller, the size of the LSP cells the for three resolution cones will be symmetric. NJ poorly identifies LSP when $a$ and $b$ are larger than $c$ and $d$ even for
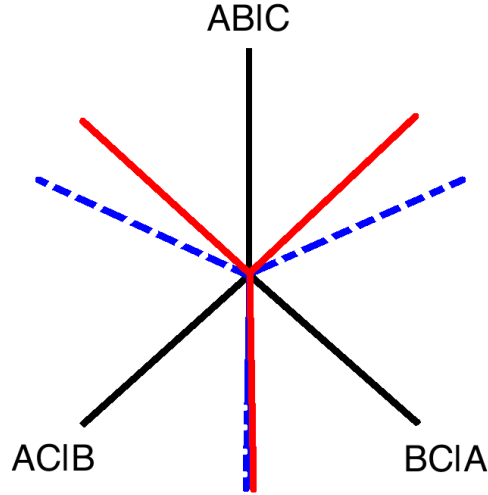
Figure 5.5: The case $a \approx b \gg c$. The cones in the fan $K_T$ are labeled with the tritomy resolutions, LSP boundaries are dashed, and UPGMA boundaries are solid.

small numbers of taxa. Unlike in the case of UPGMA, the relative size of the regions appears to be dictated not only by $a, b, c$, and $d$, but also by the topology of the subtrees with leaf sets $A, B, C$, and $D$.

**Algorithm 5.4.6** (Neighbor-Joining). • Input: a dissimilarity map $\delta \in \mathbb{R}_{\geq 0}^{n(n-1)/2}$ on $[n]$.

• Output: an unrooted binary tree $T$ and a tree metric $d$ realized by $T$.

• Initialize $[n] = \{1, 2, \ldots, n\}$, and set $d_0 = \delta$.

• For $r = 1, \ldots, n - 3$ do

  – Identify subsets $A_i, A_j$ of $[n]$ minimizing

  $$Q_r(A_i, A_j) = (n - r - 1)d_{r-1}(A_i, A_j)-$$

  $$\sum_{k=1}^{n-r+1} d_{r-1}(A_i, A_k) - \sum_{k=1}^{n-r+1} d_r - 1(A_j, A_k).$$

  – Update

  $$d_r(A_{ij}, A_k) = \frac{1}{2}d_{r-1}(A_i, A_k)+$$

  $$\frac{1}{2}d_{r-1}(A_j, A_k) - \frac{1}{2}d_{r-1}(A_i, A_j).$$

• Return: unrooted binary combinatorial tree $T$, $w : E(T) \to \mathbb{R}$ and tree metric $d_{n-3} = d$.
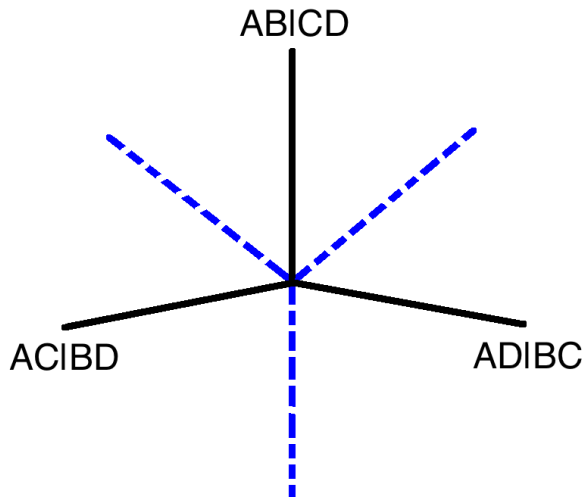
Figure 5.6: The case $a \approx b \gg c \approx d$. The cones in the fan $K_T$ are labeled with the tritomy resolutions, and LSP boundaries are dashed lines.

Applying Theorem 5.3.8, we see that when $a \approx b \gg c \approx d$, the angle between the cones $C_{T_{AC}}$ and $C_{T_{AD}}$ approaches $\pi$, while the angle measure of the LSP region for $C_{T_{AB}}$, bounded by angle bisectors between the two pairs $\{p_{AB}, p_{AC}\}$ and $\{p_{AB}, p_{AD}\}$, approaches $\pi/2$. Figure 5.6 shows this case. The fan $K_T$ is black, and the LSP boundaries are blue.

For small values of $a, b, c$, and $d$ we present computational evidence that NJ fails to identify LSP correctly. Consider the tree metrics $d_1$ and $d_2$ with topologies shown in Figures 5.7 and 5.8 and given edge weights of randomly assigned numbers between 5000 and 10000. Here $A = \{1, 2, 3, 4, 5, 6\}, B = \{7, 8, 9, 10, 11, 12\}, C = \{13\}$, and $D = \{14\}$. Using Theorem 5.3.8 we can calculate the angles between the pairs of cones in $\{C_{T_{AB}}, C_{T_{AC}}, C_{T_{AD}}\}$ and find the relative sizes of the LSP regions near the tritomies $d_1$ and $d_2$. Recall that one consequence of Theorem 5.3.8 is that the relative size of the LSP regions near $d_1$ and $d_2$ will be the same because these proportions only depend on $a, b, c$, and $d$.

Running NJ on 1,000,000 points sampled uniformly from spheres of radius 0.05 centered at $d_1$ and $d_2$ gives an empirical measure of the size of NJ regions for the three resolutions $T_{AB}, T_{AC}$, and $T_{AD}$ near the two points. Software for this experiment can be found at [27]. We compare this empirical distribution with the size of the LSP regions computed via Theorem 5.3.8 in Table 5.3. Sizes of the regions are given as percentages of the total local volume near the tritomy.

Table 5.3 shows that NJ overestimates the size of the LSP regions near $d_1$ and $d_2$ closest to the cone $C_{T_{AB}}$ and underestimates the regions near $C_{T_{AC}}$ and $C_{T_{AD}}$. Furthermore, the
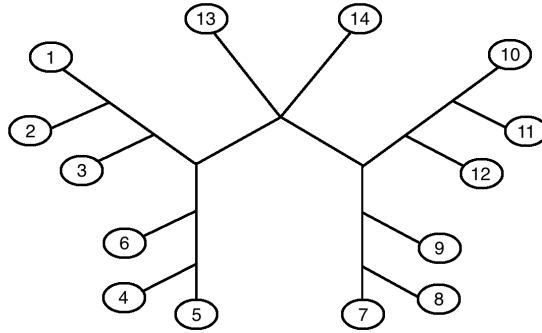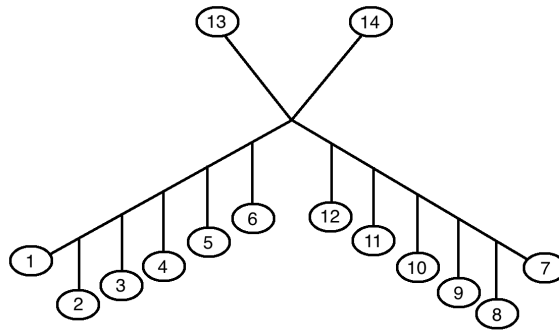
Figure 5.7: A tritomy $d_1$ on the leaf set [14].



Figure 5.8: A tritomy $d_2$ on the leaf set [14].

Table 5.3: NJ and LSP region sizes near $d_1$ and $d_2$.

| Resolution of Splits | LSP | NJ : $d_1$ | NJ: $d_2$ |
|:---:|:---:|:---:|:---:|
| $AB\|CD$ | 30.6897 | 38.1501 | 35.7037 |
| $AC\|BD$ | 34.6552 | 30.9344 | 32.1305 |
| $AD\|BC$ | 34.6552 | 30.9155 | 32.1658 |

topological structure of the subclades $A$ and $B$ influence the local size of the NJ regions. This shows that a direct analog of Theorem 5.4.5 will not exist for NJ. However, there may exist an analogous theorem for NJ when the topology of the subtrees around the polytomy is taken into account.

## 5.5 Discussion

Distance-based heuristics like UPGMA and NJ can be seen as approximating solutions to the intuitively appealing but NP-hard least-squares phylogeny problem. We compared heuristics to LSP when the true tree metric contains a tritomy. For UPGMA, our theoretical analysis shows that the success rate of the heuristic greatly depends on how balanced the sizes of the leaf sets of the underlying daughter clades are. As we discussed in Section 4.6, the decomposition of $\mathbb{R}^{n(n-1)/2}$ into NJ cones has a more complex structure. Furthermore, we have computational evidence that the geometry of NJ and LSP near a tritomy must also depend on the tree topology (Definition 1.4.2), or combinatorial branching structure of the tritomy tree. These complications have thus far prevented us from completing the same theoretical analysis for NJ near a tritomy that we performed for UPGMA. If progress is made towards finding a complete description of the family of NJ cones, we can revisit this problem.

Due to the precise form of input data used by NJ and the outputs of the algorithm, NJ is an approximation to LSP. However, Olivier Gascuel and Mike Steel showed that NJ performs a heuristic search, guided by the $Q$-criterion at each agglomeration step, that minimizes a tree-length estimate due to Yves Pauplin known as the "Balanced Minimum Evolution" (BME) criterion [39, 51]. This insight was incorporated into the selection criterion and distance re-calculation aspects of the algorithms BIONJ (due to Gascuel) [38], Weighbor (due to William Bruno, Nicholas Socci, and Aaron Halpern) [21], and FastME (due to Richard Desper and Gascuel) [29]. These algorithms take dissimilarity maps as inputs and have superior performance to NJ in terms of topological accuracy as well as better immunity to reconstruction pathologies known to biologists such as the long-branch attraction.

However, the subdivision of the input space induced by each of these improved algorithms is not polyhedral and, like the collection of Voronoi cells around a tree metric with a higher-degree polytomy, has a complicated semi-algebraic description. Any improvements to distance-based methods implied by the results in this chapter would require a fundamentally different approach, such as changing the $Q$-criterion at each step to reflect the size of the taxon groups to be joined.

# REFERENCES

[1] M. Aigner. *A Course in Enumeration.* Berlin: Springer-Verlag, 2007.

[2] D. Aldous. Stochastic models and descriptive statistics for phylogenetic trees, from Yule to today. *Statistical Science* **16** (2001): 23-34.

[3] F. Ardila and C. J. Klivans. The Bergman complex of a matroid and phylogenetic trees. *Journal of Combinatorial Theory, Series B* **96** (2006): 38-49.

[4] K. Atteson. The performance of the NJ method of phylogeny reconstruction. In: *Mathematical Hierarchies and Biology*, DIMACS Series in Discrete Mathematics and Theoretical Computer Science **37**. Editors: B. Mirkin et al. Providence, Rhode Island: American Mathematical Society (1997): 133–147.

[5] W. N. Bailey. *Generalized Hypergeometric Series.* Cambridge Tracts in Mathematics and Mathematical Physics **32**. Cambridge University Press, 1935.

[6] R. Biagioli and F. Chapoton. Supersolvable LL-lattices of binary trees. *Discrete Mathematics* **296(1)** (2005): 1–13.

[7] D. Bertsimas and J. N. Tsitsiklis. *Introduction to Linear Optimization.* Belmont, Massachusetts: Athena Scientific, 1997.

[8] L.J. Billera, S. Holmes, and K. Vogtmann. Geometry of the space of phylogenetic trees. *Advances in Applied Mathematics* **27** (2001): 733-767.

[9] L. J. Billera and A. Björner. Face numbers of polytopes and complexes. In: *Handbook of Discrete and Computational Geometry.* Editors: J. E. Goodman and J. O'Rourke. Boca Raton: CRC Press (1997): 291-310.

[10] G. Birkhoff. *Lattice Theory.* Third edition. American Mathematical Society Colloquium Publications **25**. Providence, Rhode Island: American Mathematical Society, 1967.

[11] A. Björner. Shellable and Cohen-Macaulay partially ordered sets. *Transactions of the American Mathematical Society* **260** (1980): 159-183.

[12] A. Björner. Homology and shellability of matroids and geometric lattices. In: *Matroid Applications.* Editor: N. White. Cambridge University Press (1992): 226-283.

[13] A. Björner and M. Wachs. Bruhat order of Coxeter groups and shellability. *Advances in Mathematics* **43** (1982): 87–100.

[14] A. Björner and M. Wachs. On lexicographically shellable posets. *Transactions of the American Mathematical Society* **277** (1983): 323–341.

[15] A. Björner and M. Wachs. Shellable non pure complexes and posets I. *Transactions of the American Mathematical Society* **348(4)** (1996): 1299–1327.

[16] A. Björner and M. Wachs. Shellable non pure complexes and posets II. *Transactions of the American Mathematical Society* **349(10)** (1997): 3945–3975.

[17] H. Bruggesser and P. Mani. Shellable decompositions of cells and spheres. *Mathematica Scandinavica* **29** (1971): 197-205.

[18] P. Buneman. The recovery of trees from measures of dissimilarity. In *Mathematics and the Archeological and Historical Sciences.* Editors: F. R. Hodson and D. G. Kendall. Edinburgh University Press (1971): 387-395.

[19] P. Buneman. A note on metric properties of trees. *Journal of Combinatorial Theory, Series B* **17** (1974): 48-50.

[20] G. Burde and H. Zieschang. Development of the concept of a complex. In: *History of Topology.* Editor: I. M. James. Amsterdam: North-Holland (1999): 103-110.

[21] W.M. Bruno, N.D. Socci, and A.L. Halpern. Weighted Neighbor-Joining: a likelihood-based approach to distance- based phylogeny reconstruction. *Molecular Biology and Evolution* **17** (2000): 189–197.

[22] V. Chvátal. *Linear Programming.* New York: W. H. Freeman and Company, 1983.

[23] R. Davidson and P. Hersh. A lexicographic shellability characterization of geometric lattices. *Journal of Combinatorial Theory, Series A* **123 (1)** (2014): 8–13.

[24] R. Davidson and S. Sullivant. Polyhedral combinatorics of UPGMA cones. *Advances in Applied Mathematics* **50 (2)** (2013): 327–338.

[25] R. Davidson, S. Sullivant. Supplementary materials for "Polyhedral combinatorics of UPGMA cones." [http://http://www4.ncsu.edu/ smsulli2/Pubs/UPGMACones/UPGMACones.html]

[26] R. Davidson and S. Sullivant. Distance-based phylogenetic algorithms around a polytomy. To appear in *IEEE/ACM Transactions on Computational Biology and Bioinformatics.* http://arxiv.org/abs/1307.5908

[27] R. Davidson and S. Sullivant. Software for simulating Neighbor-Joining regions near a tritomy. 2014. http://www4.ncsu.edu/∼smsulli2/Pubs/Polytomy/Polytomy.html

[28] W. Day. Computational complexity of inferring phylogenies from dissimilarity matrices. *Bulletin of Mathematical Biology* **49** (1987): 461–467.

[29] R. Desper and O. Gascuel. Theoretical foundation of the balanced minimum evolution method of phylogenetic inference and its relationship to weighted least-squares tree fitting. *Molecular Biology and Evolution* **21(3)** (2004): 587-98.

[30] A. C. Dixon. On the sum of the cubes of the coefficients in a certain expansion by the binomial theorem. *Messenger of Mathematics* **20** (1891): 7980.

[31] K. Eickmeyer, P. Huggins, L. Pachter, and R. Yoshida. On the optimality of the Neighbor-Joining algorithm. *Algorithms for Molecular Biology* **3 (5)** (2008). doi:10.1186/1748-7188-3-5

[32] K. Eickmeyer and R. Yoshida. Partitioning the sample space on five taxa for the Neighbor-Joining algorithm. http://arxiv.org/abs/math/0703081

[33] L. Euler. Solutio problematis ad geometriam situs pertinentis, Commentatio 53 indicis Enestroemiani Commentarii academiae scientiarum Petropolitanae **8** (1736): 128-140.

[34] C. Fahey, S. Hosten, N. Krieger, L. Timpe. Least squares methods for equidistant tree reconstruction. 2008. http://arxiv.org/abs/0808.3979

[35] J. Farkas. Über die Theorie der einfachen Ungleichungen. *Journal für die Reine und Angewandte Mathematik* **124** (1902): 1– 24.

[36] J. Felsenstein. *Inferring Phylogenies.* 2nd Edition. Sinauer Associates, 2003.

[37] E. Gawrilow and M. Joswig. Polymake: a framework for analyzing convex polytopes. In: *Polytopes-combinatorics and computation.* Oberwolfach (1997): 43-73. DMV Sem. 29, Birkhauser, Basel, 2000.

[38] O. Gascuel. BIONJ: an improved version of the NJ algorithm based on a simple model of sequence data. *Molecular Biology and Evolution* **14** (1997): 685-695.

[39] O. Gascuel and M. Steel. Neighbor-Joining revealed. *Molecular Biology and Evolution* **23(11)** (2006): 1997-2000.

[40] A. Garsia. Combinatorial methods in the theory of Cohen-Macaulay rings. *Advances in Mathematics* **38(3)** (1980): 229–266.

[41] J. Gill, S. Linusson, V. Moulton, and M. Steel. A regular decomposition of the edge-product space of phylogenetic trees. *Advances in Applied Mathematics* **42(2)** (2013): 158-176.

[42] A. Hatcher. *Algebraic Topology.* Cambridge: Cambridge University Press, 2002.

[43] D. Haws, T. Hodge, and R. Yoshida. Optimality of the Neighbor-Joining algorithm and faces of the balanced minimum evolution polytope. *Bulletin of Mathematical Biology* **73(11)** (2011): 2627-2648.

[44] P. Hersh, personal communication (March 2013).

[45] P. Huggins. NJBMEVolume: Software for computing volumes. 2008. http://bio.math.berkeley.edu/NJBME

[46] T.H. Jukes and C.R. Cantor. Evolution of protein molecules. *Mammalian Protein Metabolism* **3** (1969): 21-132.

[47] P. MacMahon. *Combinatory Analysis-Two Volumes in One.* New York: Chelsea Publishing, 1960.

[48] P. McMullen. The maximum number of faces of a convex polytope. *Mathematika* **17** (1970): 179-184.

[49] P. McNamara. EL-labelings, supersolvability, and 0-Hecke algebra actions on posets. *Journal of Combinatorial Theory, Series A* **101** (2003): 69-89.

[50] P. McNamara and H. Thomas. Poset edge-labelings and left modularity. *European Journal of Combinatorics* **27(1)** (2006): 101-113.

[51] Y. Pauplin. Direct calculation of a tree length using a distance matrix. *Journal of Molecular Evolution* **51** (2000): 41-47.

[52] H. Poincaré. Ir Complémeant à l'Anaylsis situs. *Rendiconti del Circolo matematico di Palermo* **13** (1899): 285-343.

[53] V. Reiner, D. Stanton, and V. Welker. The Charney-Davis quantity for certain graded posets. *Séminaire Lotharingien de Combinatoire* **50** (2003): Article B50c.

[54] I. Rival. A note on linear extensions of irreducible elements in a finite lattice. *Algebra Universalis* **6** (1976): 99–103.

[55] N. Saitou and M. Nei. The Neighbor-Joining method: a new method for reconstructing phylogenetic trees. *Molecular Biology and Evolution* **4(4)** (1987): 406-425.

[56] L. Schläfli. *Theorie der vielfachen Kontinuität*, 1852.

[57] C. Semple and M. Steel. *Phylogenetics.* Oxford University Press, 2003.

[58] J. Shareshian. On the shellability of the order complex of the subgroup lattice of a finite group. *Transactions of the American Mathematical Society* **353(7)** (2001): 2689–2703.

[59] J. M. S. Simões-Pereira. A note on the tree realizability of a distance matrix. *Journal of Combinatorial Theory* **6** (1969): 303-310.

[60] R.R. Sokal and C. Michener. A statistical method for evaluating systematic relationships. *University of Kansas Science Bulletin* **38** (1958): 1409-1438.

[61] R.R. Sokal, P.H.A. Sneath. *Numerical Taxonomy.* San Francisco: W.H. Freeman, 1963.

[62] D. Speyer and B. Sturmfels. The tropical Grassmannian. *Advances in Geometry* **4** (2004): 389-411.

[63] R. Stanley. An introduction to hyperplane arrangements. In: *Geometric Combinatorics.* Editors: E. Miller, V. Reiner, and B. Sturmfels. IAS/Park City Mathematics Series **13**. Providence, Rhode Island: American Mathematical Society, 2007.

[64] R. Stanley. Supersolvable lattices. *Algebra Universalis* **2** (1972): 197–217.

[65] R. Stanley, Finite lattices and Jordan-Hölder sets. *Algebra Universalis* **4** (1974): 361–371.

[66] R. Stanley. *Enumerative Combinatorics Volume I.* Cambridge University Press, 1997.

[67] M. Steel. The complexity of reconstructing trees from qualitative characters and subtrees. *Journal of Classification* **9** (1992): 215-233.

[68] M. Wachs. Poset topology: tools and applications. In: *Geometric Combinatorics.* Editors: E. Miller, V. Reiner, and B. Sturmfels. IAS/Park City Mathematics Series **13**. Providence, Rhode Island: American Mathematical Society, 2007.

[69] M. Wachs and J. Walker. On geometric semilattices. *Order* **2** (1986): 367–385.

[70] R. Woodroofe. An EL-labeling of the subgroup lattice. *Proceedings of the American Mathematical Society* **136(11)** (2008): 3795–3801.

[71] K. A. Zaretskii. Constructing trees from the set of distances between pendant vertices. *Uspehi Matematiceskih Nauk* **20** (1965): 90-92.

[72] G.M. Ziegler. *Lectures on Polytopes*. Graduate Texts in Mathematics **152**. New York: Springer-Verlag, 1995.