

Computational Topology

An Introduction

Herbert Edelsbrunner

John L. Harer



AMERICAN
MATHEMATICAL
SOCIETY

Computational Topology

An Introduction

Computational Topology

An Introduction

Herbert Edelsbrunner
John L. Harer



AMERICAN
MATHEMATICAL
SOCIETY

Providence, Rhode Island

“Frogzilla” cover image by Xixi Edelsbrunner.

2010 *Mathematics Subject Classification*. Primary 00-01, 52-XX, 55-XX, 57-XX, 68-XX.

For additional information and updates on this book, visit
www.ams.org/bookpages/mbk-69

Library of Congress Cataloging-in-Publication Data

Edelsbrunner, Herbert.

Computational topology : an introduction / Herbert Edelsbrunner, John L. Harer.
p. cm.

Includes bibliographical references and index.

ISBN 978-0-8218-4925-5 (alk. paper) | Softcover ISBN 978-1-4704-6769-2

1. Topology—Data processing. 2. Geometry—Data processing. 3. Computational complexity. 4. Algorithms. I. Harer, J. (John), 1952-. II. Title.

QA611.E353 2010
514—dc22

2009028121

Copying and reprinting. Individual readers of this publication, and nonprofit libraries acting for them, are permitted to make fair use of the material, such as to copy select pages for use in teaching or research. Permission is granted to quote brief passages from this publication in reviews, provided the customary acknowledgment of the source is given.

Republication, systematic copying, or multiple reproduction of any material in this publication is permitted only under license from the American Mathematical Society. Permissions to reuse portions of AMS publication content are handled by Copyright Clearance Center’s RightsLink® service. For more information, please visit: <http://www.ams.org/rightslink>.

Send requests for translation rights and licensed reprints to reprint-permission@ams.org.

Excluded from these provisions is material for which the author holds copyright. In such cases, requests for permission to reuse or reprint material should be addressed directly to the author(s). Copyright ownership is indicated on the copyright page, or on the lower right-hand corner of the first page of each article within proceedings volumes.

© 2010 by the American Mathematical Society. All rights reserved.

The American Mathematical Society retains all rights
except those granted to the United States Government.

Printed in the United States of America.

Reprinted by the American Mathematical Society 2022.

∞ The paper used in this book is acid-free and falls within the guidelines
established to ensure permanence and durability.
Visit the AMS home page at <http://www.ams.org/>
13 12 11 10 9 8 7 6 5 27 26 25 24 23 22

To our parents,

Herbert and Berta Edelsbrunner,
Chuck and Prinny Harer

Contents

Preface	xi
A Computational Geometric Topology	1
I Graphs	3
I.1 Connected Components	3
I.2 Curves in the Plane	9
I.3 Knots and Links	13
I.4 Planar Graphs	18
Exercises	24
II Surfaces	27
II.1 2-dimensional Manifolds	27
II.2 Searching a Triangulation	33
II.3 Self-intersections	37
II.4 Surface Simplification	42
Exercises	47
III Complexes	51
III.1 Simplicial Complexes	51
III.2 Convex Set Systems	57
III.3 Delaunay Complexes	63
III.4 Alpha Complexes	68
Exercises	74

B Computational Algebraic Topology	77
IV Homology	79
IV.1 Homology Groups	79
IV.2 Matrix Reduction	85
IV.3 Relative Homology	90
IV.4 Exact Sequences	95
Exercises	101
V Duality	103
V.1 Cohomology	103
V.2 Poincaré Duality	108
V.3 Intersection Theory	114
V.4 Alexander Duality	118
Exercises	123
VI Morse Functions	125
VI.1 Generic Smooth Functions	125
VI.2 Transversality	130
VI.3 Piecewise Linear Functions	135
VI.4 Reeb Graphs	140
Exercises	145

C Computational Persistent Topology	147
VII Persistence	149
VII.1 Persistent Homology	149
VII.2 Efficient Implementations	156
VII.3 Extended Persistence	161
VII.4 Spectral Sequences	166
Exercises	171
VIII Stability	175
VIII.1 1-parameter Families	175
VIII.2 Stability Theorems	180
VIII.3 Length of a Curve	185
VIII.4 Bipartite Graph Matching	191
Exercises	197
IX Applications	199
IX.1 Measures for Gene Expression Data	199
IX.2 Elevation for Protein Docking	206
IX.3 Persistence for Image Segmentation	213
IX.4 Homology for Root Architectures	218
Exercises	224
References	227
Index	235

Preface

The last ten years have witnessed the fact that geometry, topology, and algorithms form a potent mix of disciplines with many applications inside and outside academia. We aim at bringing these developments to a larger audience. This book has been written to be taught, and it is based on notes developed during courses delivered at Duke University and at the Berlin Mathematical School, primarily to students of computer science and mathematics. The organization into chapters, sections, and exercises reflects the teaching style we practice. Each chapter develops a major topic and provides material for about two weeks. The chapters are divided into sections, each a lecture of one and a quarter hours. An interesting challenge is the mixed background of the audience. How do we teach topology to students with a limited background in mathematics, and how do we convey algorithms to students with a limited background in computer science? Assuming no prior knowledge and appealing to the intelligence of the listener are good first steps. Motivating the material by relating it to situations in different walks of life is helpful in building up intuition that can cut through otherwise necessary formalism. Exposing central ideas with simple means helps, and so does minimizing the necessary amount of technical detail.

The material in this book is a combination of topics in geometry, topology, and algorithms. Far from getting diluted, we find that the fields benefit from each other. Geometry gives a concrete face to topological structures, and algorithms offer a means to construct them at a level of complexity that passes the threshold necessary for practical applications. As always, algorithms have to be fast because time is the one fundamental resource humankind has not yet learned to manipulate for its selfish purposes. Beyond these obvious relationships, there is a symbiotic affinity between algorithms and the algebra used to capture topological information. It is telling that both fields trace their names back to the writing of the same Persian mathematician, al-Khwarizmi, working in Baghdad during the ninth century after Christ. Besides living in the triangle spanned by geometry, topology, and algorithms, we find it useful to contemplate the place of the material in the tension between extremes such as local vs. global, discrete vs. continuous, abstract vs. concrete, and intrinsic vs. extrinsic. Global insights are often obtained by a meaningful integration of local information. This is how we proceed in many fields, taking on bigger challenges after mastering the small ones. But small things are big from up close, and big things are small from afar. Indeed, the question of scale lurking behind this thought is the

driving force for much of the development described in this book. The dichotomy between discrete and continuous structures is driven by opposing goals: machine computation and human understanding. The tension between the abstract and the concrete as well as between the intrinsic and the extrinsic has everything to do with the human approach to knowledge. An example close to home is the step from geometry to topology in which we remove the burdens of size to focus on the phenomenon of connectivity. The more abstract the context the more general the insight. Now, generality is good, but it is not a substitute for the concrete steps that have to be taken to build bridges to applications. Zooming in and out of generality leads to unifying viewpoints and suggests meaningful integrations where they exist.

While these thoughts have certainly influenced us in the selection of the material and in its presentation, there is a long way to the concrete instantiation we call this book. It consists of three parts and nine chapters. Part A is a gentle introduction to topological thought. Discussing graphs in Chapter I, surfaces in Chapter II, and complexes in Chapter III, we gradually build up topological sophistication, always in combination with geometric and algorithmic ideas. Part B presents classical material from topology. We focus on what we deem useful and efficiently computable. The material on homology in Chapter IV and duality in Chapter V is exclusively algebraic. In the discussion of Morse theory in Chapter VI, we build a bridge to differential concepts in topology. Part C is novel and the reason for why we wrote this book. The main new concept is persistence, introduced in Chapter VII, and its stability, discussed in Chapter VIII. Finally, we discuss applications in Chapter IX.

In a project like writing this book, there are many who contribute, directly or indirectly. We want to thank all, but we don't know where to begin. Above all, we thank our colleagues in academia and industry, our students, and our postdoctoral fellows for their ideas, criticism, and encouragement, and most of all for the sense of purpose they instilled. We thank Duke University and IST Austria for providing the facilities and intellectual environment that allowed us to engage in the line of research leading to this book. We thank the computer science and the mathematics departments at Duke University and the Berlin Mathematical School for the opportunity to teach computational topology to their students. These courses provided the motivation to develop the notes that turned into this book. We are grateful to the funding agencies for nurturing the research that led to this book. The National Science Foundation and the National Institute of Health generously supported our collaborations with biochemists and biologists. Most of all, we thank our program manager at the Defense Advanced Research Projects Agency, Benjamin Mann, for his continued support and his enthusiasm for our research. Last but not least, we thank Ina Mette for believing in this project and the staff at the American Mathematical Society for making the steps toward the final product an enjoyable experience.

Herbert Edelsbrunner and John L. Harer
Durham, North Carolina, 2009

Part A

Computational Geometric Topology

Chapter I

Graphs

In topology we think of a graph as a 1-dimensional geometric object, vertices being points and edges being curves connecting these points in pairs. This view is different from but compatible with the interpretation of a graph common in discrete mathematics where the vertices are abstract elements and the edges are pairs of these elements. In more than one way, this book lives in the tension between the discrete and the continuous, and graphs are just one example of this phenomenon. We begin with the discussion of an intrinsic property, namely whether a graph is connected or not. Indeed, this does not depend on where we draw the graph, on paper or in the air. Following are extrinsic considerations about curves and graphs in the plane and in 3-dimensional space. While the extrinsic questions are natural to most people, the mathematician usually favors the intrinsic point of view since it tends to lead to more fundamental insights of more general validity.

I.1 Connected Components

A theme that goes through this entire book is the exchange between discrete and continuous models of reality. In this first section, we compare the notion of connectedness in discrete graphs and continuous spaces.

Simple graphs. An abstract *graph* is a pair $G = (V, E)$ consisting of a set of *vertices*, V , and a set of *edges*, E , each a pair of vertices. We draw the vertices as points or little circles and the edges as line segments or curves connecting the points. The graph is *simple* if the edge set is a subset of the set of unordered pairs, $E \subseteq \binom{V}{2}$, which means that no two edges connect the same two vertices and no edge joins a vertex to itself. For $n = \text{card } V$ vertices, the number of edges is $m = \text{card } E \leq \binom{n}{2}$. Every simple graph with n vertices is a subgraph of the *complete graph*, K_n , that contains an edge for every pair of vertices; see Figure I.1.

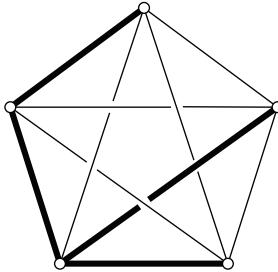


Figure I.1: The complete graph with five vertices, K_5 . It has ten edges which form five crossings if drawn as sides and diagonals of a convex pentagon. The four thick edges connect the same five vertices and form a spanning tree of the complete graph.

In a simple graph, a *path* between vertices u and v can be described by a sequence of vertices, $u = u_0, u_1, u_2, \dots, u_k = v$, with an edge between u_i and u_{i+1} for each $0 \leq i \leq k-1$. The *length* of this path is its number of edges, k . Vertices can repeat, allowing the path to cross itself or backtrack. The path is *simple* if the vertices in the sequence are distinct, that is, $u_i \neq u_j$ whenever $i \neq j$.

DEFINITION. A simple graph is *connected* if there is a path between every pair of vertices.

A (*connected*) *component* is a maximal subgraph that be connected. The smallest connected graphs are the *trees*, which are characterized by having a unique simple path between every pair of vertices. Removing any one edge disconnects the tree. A *spanning tree* of $G = (V, E)$ is a tree (V, T) with $T \subseteq E$; see Figure I.1. It has the same vertex set as the graph and uses a minimal set of edges necessary to be connected. This requires that the graph be connected to begin with. Indeed, a graph is connected iff it has a spanning tree. An alternative characterization of a connected graph can then be based on the impossibility to cut it in two.

DEFINITION. A *separation* is a non-trivial partition of the vertices; that is, $V = U \dot{\cup} W$ with $U, W \neq \emptyset$, such that no edge connects a vertex in U with a vertex in W . A simple graph is *connected* if it has no separation.

Topological spaces. We now switch to a continuous model of reality, the topological space. There are similarities to graphs, in particular if our interest is limited to questions of connectedness. Starting with a point set, we consider a topology, which is a way to define which points are near without specifying how near they are from each other. Think of it as an abstraction of Euclidean space in which neighborhoods are open balls around points. Concretely, a *topology* on a point set \mathbb{X} is a collection \mathcal{U} of subsets of \mathbb{X} , called *open sets*, such that

- (i) \mathbb{X} is open and the empty set \emptyset is open;
- (ii) if U_1 and U_2 are open, then $U_1 \cap U_2$ is open;
- (iii) if U_i is open for all i in some possibly infinite, possibly uncountable, index set, then the union of all the U_i is open.

The pair $(\mathbb{X}, \mathcal{U})$ is called a *topological space*, but we will usually tacitly assume that \mathcal{U} is understood and refer to \mathbb{X} as a topological space. Since we can repeat the pairwise intersection, condition (ii) is equivalent to requiring that common intersections of finitely many open sets be open.

To build interesting topologies, we start with some initial notion of which sets might be open and then form appropriate combinations of these until the three conditions are satisfied. A *basis* of a topology on a point set \mathbb{X} is a collection \mathcal{B} of subsets of \mathbb{X} , called *basis elements*, such that each $x \in \mathbb{X}$ is contained in at least one $B \in \mathcal{B}$ and $x \in B_1 \cap B_2$ implies there is a third basis element with $x \in B_3 \subseteq B_1 \cap B_2$. The topology \mathcal{U} generated by \mathcal{B} consists of all sets $U \subseteq \mathbb{X}$ for which $x \in U$ implies there is a basis element $x \in B \subseteq U$. This topology can be constructed explicitly by taking all possible unions of all possible finite intersections of basis sets. As an example consider the real line, $\mathbb{X} = \mathbb{R}$, and let \mathcal{B} be the collection of open intervals. This defines the usual topology of the real line. Note that the intersection of the intervals $(-\frac{1}{k}, +\frac{1}{k})$, for the infinitely many integers $k \geq 1$, is the point 0. This is not an open set, which illustrates the need for the restriction to finite intersections.

We often encounter sets inside other sets, $\mathbb{Y} \subseteq \mathbb{X}$, and in these cases we can borrow the topology of the latter for the former. Specifically, if \mathcal{U} is a topology of \mathbb{X} , then the collection of sets $\mathbb{Y} \cap U$, for $U \in \mathcal{U}$, is the *subspace topology* of \mathbb{Y} . As an example consider the closed interval $[0, 1] \subseteq \mathbb{R}$. We have seen that the open intervals form a basis for a topology of the real line. The intersections of open intervals with $[0, 1]$ form the basis of the subspace topology of the closed interval. Note that the interval $(1/2, 1]$ is considered an open set in $[0, 1]$, but it is not open when considered as a set in \mathbb{R} .

Continuity, paths, and connectedness. A function from one topological space to another is *continuous* if the preimage of every open set is open. This is derived from the concept of continuity familiar from calculus; for example the function $f : \mathbb{R} \rightarrow \mathbb{R}$ that maps $(-\infty, 0]$ to 0 and $(0, \infty)$ to 1 is not continuous because for any $0 < \varepsilon < 1$, $(-\varepsilon, \varepsilon)$ is open, but $f^{-1}((-\varepsilon, \varepsilon))$ is not.

A *path* is a continuous function from the unit interval, $\gamma : [0, 1] \rightarrow \mathbb{X}$. It *connects* the point $\gamma(0)$ to the point $\gamma(1)$ in \mathbb{X} . Similar to paths in graphs, we allow for self-intersections, that is, values $s \neq t$ in the unit interval for which $\gamma(s) = \gamma(t)$. If there are no self-intersections, then the function is injective and the path is *simple*. Now we are ready to adapt our first definition of connectedness to topological spaces.

DEFINITION. A topological space is *path-connected* if every pair of points is connected by a path.

There is also a counterpart of our second definition of connectedness. We formulate it using open sets, and there is an equivalent formulation in terms of *closed sets* which, by definition, are complements of open sets.

DEFINITION. A *separation* of a topological space \mathbb{X} is a partition $\mathbb{X} = U \dot{\cup} W$ into two non-empty, open subsets. A topological space is *connected* if it has no separation.

It turns out connectedness is strictly weaker than path-connectedness, although for most spaces we will encounter they are the same. An example of a space that is connected but not path-connected is the comb with a single tooth deleted. It is constructed by gluing vertical teeth to a horizontal bar and finally deleting the interior of the last tooth: taking the union of $[0, 1] \times 0$, $0 \times [0, 1]$, $\frac{1}{k} \times [0, 1]$, for all positive integers k , we finally delete $0 \times (0, 1)$. To construct a topology, we take the collection of open disks as the basis of a topology on \mathbb{R}^2 and we use the subspace topology for the comb. This space is connected because it is the union of a path-connected set and a limit point. It is not path-connected because no path from anywhere else can reach 0×1 .

Disjoint set systems. We return to graphs and consider the algorithmic question of deciding connectedness. There are many approaches, and we present a solution based on maintaining a disjoint set system. This particular algorithm has various other applications, some of which will be discussed in later chapters of this book. Using the integers from 1 to n as the names of the vertices, we store each component of the graph as a subset of $[n] = \{1, 2, \dots, n\}$. Adding the edges one at a time and maintaining the system of sets representing the components, we find that the graph is connected iff in the end there is only one set left, namely $[n]$. Formulated as an abstract data type, we have two operations manipulating the system:

FIND(i): return the name of the set that contains i ;

UNION(i, j): assuming i and j belong to different sets in the system, replace the two sets by their union.

We need the find operation to test whether i and j indeed belong to different sets. Each successful union operation reduces the number of sets in the system by one. Starting with n singleton sets, it therefore takes $n - 1$ union operations to get to a single set. Since trees connecting the n vertices can be generated this way, we thus have a proof that every tree with n vertices has $m = n - 1$ edges.

A standard data structure implementing a disjoint set system stores each set as a tree embedded in a linear array, $V[1..n]$. Each node in the tree is equipped with a pointer to its *parent*, except for the *root* which has no parent; see Figure I.2. Who is parent of whom is not important as long as the vertices are connected in a single tree. We implement the find operation by traversing the tree upward until we reach the root, reporting the root as the name of the set.

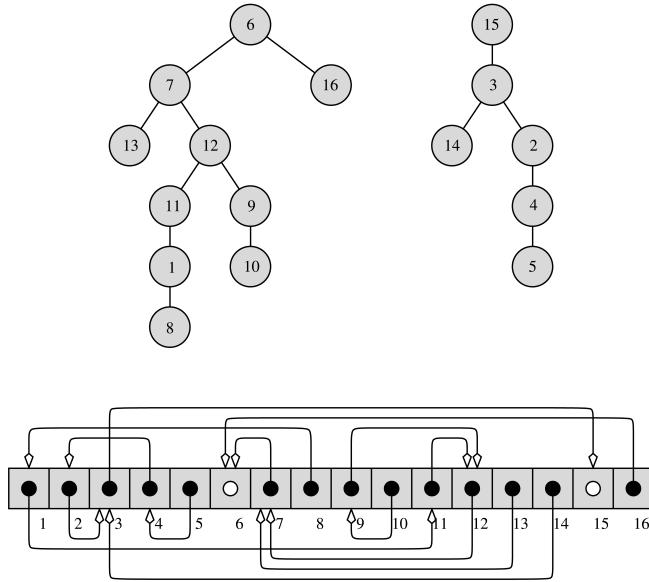


Figure I.2: Top: two trees representing two disjoint sets. Bottom: storing the two trees in a linear array using an arbitrary ordering of the nodes.

```
int FIND( $i$ )
    if  $V[i].parent \neq \text{null}$  then return FIND( $V[i].parent$ )
        else return  $i$ 
    endif.
```

If i is not the root, then we find the root recursively and finally return it. Otherwise, we return i as the root. We implement the union operation by linking one root to the other.

```
void UNION( $i, j$ )
 $x = \text{FIND}(i); y = \text{FIND}(j);$ 
if  $x \neq y$  then  $V[x].parent = y$  endif.
```

After making sure that the two sets are different, we assign one root as the parent of the other.

Improving the running time. The above implementation is not very efficient, in particular if we have long paths that are repeatedly traversed. To prevent this from happening, we always link the smaller to the larger tree.

```

void UNION( $i, j$ )
   $x = \text{FIND}(i); y = \text{FIND}(j);$ 
  if  $x \neq y$  then
    if  $V[x].size > V[y].size$  then  $x \leftrightarrow y$  endif;
     $V[x].parent = y$ 
  endif.

```

Now a tree of k nodes cannot have paths longer than $\log_2 k$ edges since the size of the subtree grows by at least a factor of two each time we pass to the parent. To further improve the efficiency, we compress paths whenever we traverse them. Here it is convenient to assume that roots are identified by pointing to themselves.

```

int FIND( $i$ )
  if  $V[i].parent \neq i$  then
    return  $V[i].parent = \text{FIND}(V[i].parent)$ 
  endif;
  return  $i$ .

```

If i is not the root, then the function recursively finds the root, makes the root the parent of i , reports the root, and exits. Otherwise, the function reports i as the root and exits.

In analyzing the algorithm, we are interested in the running time for executing a sequence of m union and find operations. Finding tight bounds turns out to be a difficult problem, and we limit ourselves to stating the result. Specifically, if n is the number of vertices, then the running time is never more than some constant times $m\alpha(n)$, where $\alpha(n)$ is the notoriously slow growing inverse of the Ackermann function. Eventually, $\alpha(n)$ goes to infinity, but to reach even beyond five, we need an astronomically large number of vertices, more than the estimated number of electrons in our Universe. In other words, for all practical purposes the algorithm takes constant average time per operation, but theoretically this is not a true statement.

Bibliographic notes. Graphs are ubiquitous objects and appear in most disciplines. Within mathematics, the theory of graphs is considered part of combinatorics. There are many good books on the subject, including the one by Tutte [142]. The basic topological notions of connectedness are treated in books on point-set or general topology, including the text by Munkres [115]. The computational problem of maintaining a system of disjoint sets, often referred to as the union-find problem, is a classic topic in the field of algorithms. Solutions to it are known as union-find data structures, and the most efficient of all is the up-tree representation maintained through weighted union and path-compression as explained in this section. A complete description of the non-trivial analysis of the algorithm can be found in the text by Tarjan [140].

I.2 Curves in the Plane

In the previous section, we used paths to merge points into connected components. To capture aspects of connectivity that go beyond components, we need different maps.

Closed curves. We distinguish primarily between two kinds of (connected) curves, *paths* and *closed curves*. As defined in the previous section, paths are continuous maps from $[0, 1]$ to \mathbb{X} . Sometimes, a closed curve is defined as a path in which 0 and 1 map to the same point. Usually, we will define a closed curve to be a map from the unit circle, $\gamma : \mathbb{S}^1 \rightarrow \mathbb{X}$, where $\mathbb{S}^1 = \{x \in \mathbb{R}^2 \mid \|x\| = 1\}$. This second version emphasizes the important fact that paths and closed curves capture different properties of topological spaces, since the interval and the circle are different topological spaces. To make this precise, we call two topological spaces *homeomorphic* or *topologically equivalent* if there exists a continuous bijection from one space to the other whose inverse is also continuous. A map with these properties is called a *homeomorphism*. Notice that a homeomorphism between two spaces gives a bijection between their open sets. The unit interval and the unit circle are not homeomorphic. Indeed, removing the midpoint decomposes the interval into two components while removing any point leaves the circle connected. This contradicts the existence of a bijection that is continuous in both directions.

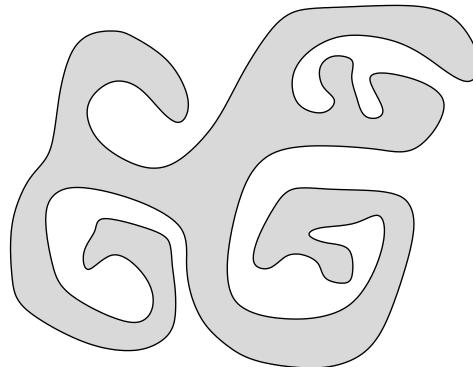


Figure I.3: The shaded inside and the white outside of a simple closed curve in the plane.

Considering maps into the Euclidean plane, $\mathbb{X} = \mathbb{R}^2$, it makes sense to distinguish curves with and without self-intersections. A *simple closed curve* is a curve without self-intersections, that is, a continuous injection $\gamma : \mathbb{S}^1 \rightarrow \mathbb{R}^2$. Interestingly, every such curve decomposes the plane into two pieces, one on each side of the curve, as in Figure I.3.

JORDAN CURVE THEOREM. Removing the image of a simple closed curve from \mathbb{R}^2 leaves two connected components, the bounded *inside* and the unbounded *outside*.

The inside together with the image of the curve is in fact homeomorphic to a closed disk.

This may seem obvious, but proving it is challenged by the generality of the claim which is formulated for all and not just smooth or piecewise linear simple closed curves. There are reasons to believe that there is no simple proof for such a general claim. The fact that the inside together with the curve is homeomorphic to the closed disk, $\mathbb{B}^2 = \{x \in \mathbb{R}^2 \mid \|x\| \leq 1\}$, is known as the Schönflies Theorem. The Jordan Curve Theorem remains valid if we replace the plane by the sphere, $\mathbb{S}^2 = \{x \in \mathbb{R}^3 \mid \|x\| = 1\}$, but not if we replace it by the torus. The Schönflies Theorem as stated is false in dimension higher than two.

Parity algorithm. Given a simple closed curve in the plane, a fundamental computational question asks whether a given *query point* $x \in \mathbb{R}^2$ lies inside, on, or outside the curve. To write an algorithm answering this question, we assume a finite approximation of the curve. For example, we may specify γ at a finite number of points and interpolate linearly between them. The result is a *closed polygon*; see Figure I.4. It is *simple* if it is a closed curve itself. To decide whether the point x

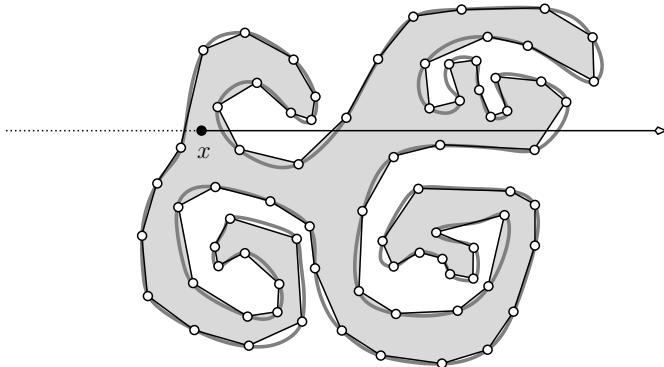


Figure I.4: Approximation of the simple closed curve in Figure I.3 by a simple closed polygon. The point x lies inside the polygon and the half-line crosses the polygon an odd number of times.

lies inside such a simple closed polygon, we draw a half-line emanating from x and count how often it crosses the polygon. Assuming x does not lie on the polygon, it lies inside if the number of crossings is odd and outside if that number is even, hence the name Parity Algorithm. In the implementation of this idea, we let $x = (x_1, x_2)$ be the query point and $a = (a_1, a_2)$, $b = (b_1, b_2)$ the endpoints of an edge of the polygon. We assume the generic case in which no three points are collinear and no two lie on a common vertical or horizontal line. To simplify the code, we choose the horizontal half-line leaving x toward the right and we assume that a is below b , that is, $a_2 < b_2$. We first make sure that the entire horizontal line crosses the edge, which we do by testing $a_2 < x_2 < b_2$. If it does, then we test whether the

crossing lies to the left or the right of the query point. To this end we compute the determinant of the matrix

$$\Delta(x, a, b) = \begin{bmatrix} 1 & x_1 & x_2 \\ 1 & a_1 & a_2 \\ 1 & b_1 & b_2 \end{bmatrix},$$

which is positive iff the sequence of points x, a, b forms a left-turn. To see this, we verify the claim for $x = (0, 0)$, $a = (1, 0)$, $b = (0, 1)$ and then notice that the sign of the determinant switches exactly when the three points become collinear. We use this fact to decide whether the half-line crosses the edge:

```
boolean DOESCROSS( $x, a, b$ )
  if not  $a_2 < x_2 < b_2$  then return FALSE endif;
  return  $\det \Delta(x, a, b) > 0$ .
```

Now we run this test for all edges and this way count the crossings. The trouble with this implementation is the non-generic cases. We finesse them using two infinitesimally small, positive numbers $0 < \varepsilon_1 \ll \varepsilon_2$ and substituting $x' = (x_1 + \varepsilon_1, x_2 + \varepsilon_2)$ for x . A generic case for x is generic for x' , and we get the same decision for both points. A non-generic case for x is generic for x' , and we use the decision for x' .

Polygon triangulation. Sometimes it is useful to have a more structured representation of the inside of the polygon, for example for navigation to find the exit out of a maze. The most common such structural representation is a *triangulation* which is a decomposition into triangles. Here we require that the triangles use the vertices of the polygon but do not introduce new ones. Furthermore, they use the edges of the polygon together with *diagonals*, which are new edges that connect non-adjacent vertices of the polygon. The diagonals are required to pass through the inside and not cross any other diagonals and any polygon edges; see Figure I.5.

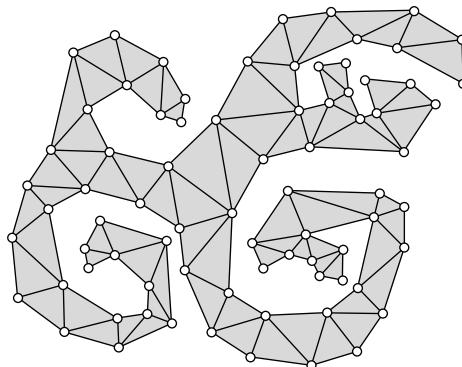


Figure I.5: A triangulation of the polygon in Figure I.4. Each diagonal passes from one side of the inside to the other.

To prove that a triangulation always exists, we just need to show that there is at least one diagonal, unless the number of edges in the polygon is $n = 3$. Indeed, we may consider the leftmost vertex, b , of the polygon. Either we can connect its two neighbors, a and c , or we can connect b to the leftmost vertex u that lies inside the triangular region abc . Drawing this diagonal decomposes the n -gon into two, an n_1 -gon and an n_2 -gon. We have $n_1 + n_2 = n + 2$, and since both are at least three, we also have $n_1, n_2 < n$. We can thus use induction to complete the proof. The same inductive argument shows that there are $n - 3$ diagonals and $n - 2$ triangles, no matter how we triangulate. This is suggestive. Indeed, we can think of the triangles as the nodes and the diagonals as the arcs of a tree. Since every tree with $n - 2 \geq 2$ nodes has at least one leaf, that is, a node with only one neighbor, every triangulation has an *ear*, that is, a triangle formed by one diagonal and two polygon edges. Incidentally, this is another property that does not generalize to tetrahedral decompositions in \mathbb{R}^3 .

Winding number. We return to a general, not necessarily simple, closed curve $\gamma : \mathbb{S}^1 \rightarrow \mathbb{R}^2$. Let x be a point not in the image of the curve. Suppose we traverse γ and view the moving point from x . Specifically, we let s go once around the circle and observe the unit vector $(\gamma(s) - x)/\|\gamma(s) - x\|$ rotate about the origin. When the vector completes a full turn, we count $+1$ or -1 depending on whether this turn is counterclockwise or clockwise. The sum of these numbers is the *winding number* of γ and x , denoted as $W(\gamma, x)$. It is necessarily an integer and gives the net number of counterclockwise turns we observe. If γ is simple, then the points inside the curve all have the same winding number, -1 or $+1$. To reduce this to one case, we may reorient the curve, e.g. by reflecting the unit circle along the horizontal coordinate axis, and get

$$W(\gamma, x) = \begin{cases} +1 & \text{if } x \text{ is inside;} \\ 0 & \text{if } x \text{ is outside.} \end{cases}$$

However, for non-simple curves we can get winding numbers of absolute value larger than one; see Figure I.6. Suppose we move x in the plane. As long as it does not cross the curve, the winding number does not change. Crossing the curve changes the winding number, namely by -1 if we cross from left to right and by $+1$ if we cross from right to left. But this implies that at least two regions in the decomposition defined by γ have their boundary arcs consistently oriented. Specifically, the neighbors of a region with locally maximum winding number all have winding number one less, so the region lies to the left of all its boundary arcs. Similarly, a region with locally minimal winding number lies to the right of all its boundary arcs.

Bibliographic notes. The Jordan Curve Theorem is well known also beyond topology, in part because it seems so obvious but at the same time is difficult to prove. The theorem is named after Camille Jordan, who was the first to present a proof in 1882. This proof was later found to be incorrect, and the first satisfactory proof is due to Veblen [145]. We refer to [151] for a more recent discussion of the result.

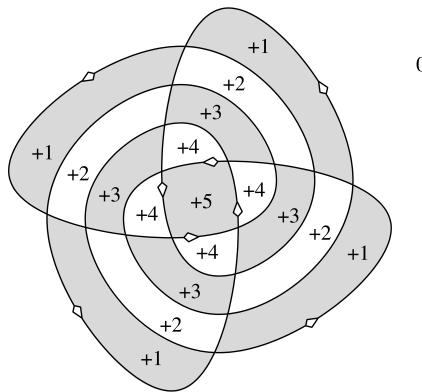


Figure I.6: An oriented non-simple closed curve with regions distinguished by the winding number of their points.

The difficulties encountered in the implementation of the Parity Algorithm for locating a point inside or outside a closed polygon have been voiced in [71]. A provably correct implementation can be achieved with exact arithmetic and symbolic perturbation as described in [62]. Triangulations of simple closed polygons in the plane have been studied in computational geometry. Constructing such a triangulation in time proportional to the number of vertices seems rather difficult and the algorithm by Chazelle [32] that achieves this feat is not recommended for implementation.

I.3 Knots and Links

In this section, we study closed curves in 3-dimensional Euclidean space and questions how they relate to each other and to themselves.

Knots. A closed curve embedded in \mathbb{R}^3 does not decompose the space, but it can be tangled up in inescapable ways. The field of mathematics that studies such tangles is knot theory. Its prime subject is a *knot* which is an *embedding* $\kappa : \mathbb{S}^1 \rightarrow \mathbb{R}^3$, that is, an injective, continuous function that is a homeomorphism onto its image. It turns out that any injective, continuous function from \mathbb{S}^1 to \mathbb{R}^3 is an embedding, but this is not true for general domains. Another knot is *equivalent* to κ if it can be continuously deformed into κ without crossing itself during this process. Equivalent knots are considered the same. The simplest knot is the *unknot*, also known as the *trivial knot*, which can be deformed to a geometric circle in \mathbb{R}^3 . Two other and only slightly tangled up knots are the *trefoil knot* and the *figure-eight knot*, both illustrated in Figure I.7. A subtlety in the definition of equivalence is that deformations in which knotted parts disappear in the limit are not allowed. It is therefore useful to think of knots as curves with small but positive thickness, similar to shoelaces and ropes.

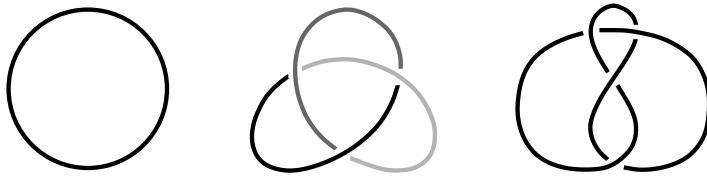


Figure I.7: From left to right: the unknot, the trefoil knot, and the figure-eight knot. The trefoil knot is tricolored.

Reidemeister moves. Let us follow a deformation of a knot by drawing its projections to a plane, keeping track of the underpasses and overpasses at crossings. We are primarily interested in generic projections defined by the absence of any violations to injectivity, other than a discrete collection of double points where the curve crosses itself in the plane. In a generic deformation, we observe three types of non-generic projections that transition between generic projections, which are illustrated in Figure I.8. It is plausible and also true that any two generic projections of the same knot can be transformed into each other by Reidemeister moves.

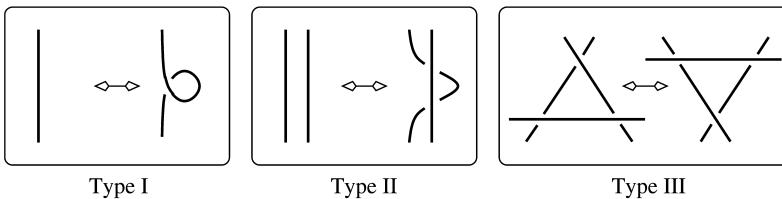


Figure I.8: The three types of Reidemeister moves.

It seems clear that the trefoil knot is not equivalent to the unknot, and there is indeed an elementary proof using Reidemeister moves. Call a piece of the knot from one underpass to the next a *strand*. A *tricoloring* of a generic projection colors each strand with one of three colors such that

- (i) at each crossing either three colors or only one color come together;
- (ii) at least two colors are used.

Figure I.7 shows that the standard projection of the trefoil knot is tricolorable. A useful property of Reidemeister moves is that they preserve tricolorability; that is, the projection before the move is tricolorable iff the projection after the move is tricolorable.

TYPE I. Going from left to right in Figure I.8, we use the same one color, and going from right to left, we observe that we have only one color coming together at the crossing.

TYPE II. From left to right we have two possibilities, either using only one color or going from two to three colors. The reverse direction is symmetric.

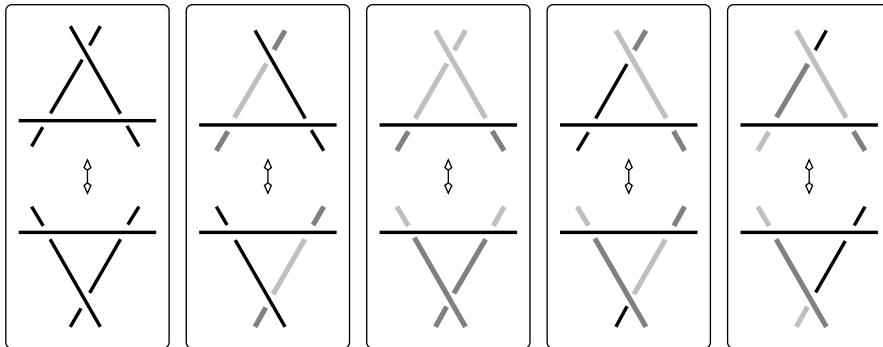


Figure I.9: The different cases in the proof that the Type III Reidemeister move preserves tricolorability. In each case there is only one new strand whose color can be chosen anew.

TYPE III. There are five cases to be checked, all shown in Figure I.9.

The trefoil knot is tricolorable and the unknot is not tricolorable. It follows that the two are not equivalent. It is not difficult to see that the figure-eight knot is not tricolorable. This implies that the trefoil knot and the figure-eight knot are different, but the method is not powerful enough to distinguish the figure-eight from the unknot.

Links. A *link* is a collection of two or more disjoint knots. Equivalence between links is defined the same way as between knots, and Reidemeister moves again suffice to go from one generic projection to another. A disjoint plane *splits* a link if there are knots on both sides. A link is *splittable* if an equivalent link has a splitting plane. The *unlink* or *trivial link* of size two consists of two unknots that can be split, like the two circles in Figure I.10 on the left. The easiest non-splittable link consisting of two unknots is the *Hopf link*, which is shown in Figure I.10 in the middle. We can again use tricolorability to prove that the Hopf link is different

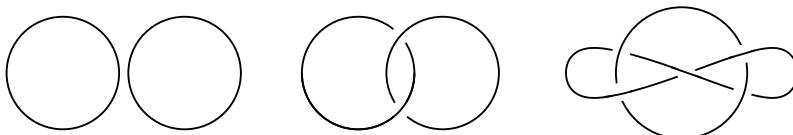


Figure I.10: From left to right: the unlink, the Hopf link, and the Whitehead link.

from the unlink. Alternatively, we may count the crossings between the two knots, κ and λ , counting with a sign. Specifically, we orient each knot arbitrarily and we look at each crossing locally. If the underpass goes from the left of the overpass to its right, then we count $+1$, and otherwise we count -1 . Letting x be a crossing and $\text{sign}(x)$ be plus or minus 1 as explained, the *linking number* is half the sum of

these numbers over all crossings:

$$Lk(\kappa, \lambda) = \frac{1}{2} \sum_x \text{sign}(x).$$

Changing the orientation of one knot but not the other has the effect of reversing the sign of the linking number. Clearly, Reidemeister moves do not affect the linking number. Since the linking number of the unlink is zero and that of the Hopf link is plus or minus 1, we have another proof that the two links are different. An easy link that is not splittable but has vanishing linking number is the Whitehead link in Figure I.10. It consists of two unknots but cannot be tricolored, which implies that it is not splittable.

Writhing number. Next we introduce a number that measures how contorted the curve is in space. Let $\kappa : \mathbb{R}^1 \rightarrow \mathbb{R}^3$ be a knot and assume that it is smooth and its tangent vector $\dot{\kappa}(s)$ is non-zero for every s . Projecting along a direction $u \in \mathbb{S}^2$, we get a closed curve in the plane. Assuming the projection is generic, we distinguish underpass from overpasses and count each crossing plus or minus 1 time, as before. However, different from the case of the linking number, we count crossings that the curve makes with itself and we do not divide by two. The sum of these numbers is the *directional writhing number*, $DWr(\kappa, u)$. The *writhing number* is the average over all directions. This is the integral of the directional writhing number over all directions divided by the area of the unit sphere:

$$Wr(\kappa) = \frac{1}{4\pi} \int_{u \in \mathbb{S}^2} DWr(\kappa, u) du.$$

The directions with non-generic projections form only a measure zero subset of the sphere. We therefore make no mistake when we average only over all generic projections. In contrast to the linking number, the writhing number is not necessarily an integer and it depends on the exact shape of the curve. Besides the shape it also captures topological information, as we will see shortly.

A good motivation for studying the writhing number comes from molecular biology and, more specifically, the shape of DNA within the cell. Modeling its double-helix structure with a constant width ribbon, we are interested in the writhing number of the center axis, κ . The boundary of the ribbon consists of two closed curves. We need only one, $\lambda : \mathbb{S}^1 \rightarrow \mathbb{R}^3$. In the case of DNA, λ twists and turns around κ . Intuitively, the *twisting number* is the average motion of λ relative to κ . To formalize this idea, we assume that the center axis and the boundary curve are one unit of length apart and parametrized such that $\lambda(s) - \kappa(s)$ has unit length and is normal to the center axis. We construct a frame of mutually orthogonal unit vectors consisting of the tangent vector at s , $T(s) = \dot{\kappa}(s)/\|\dot{\kappa}(s)\|$, the normal vector connecting the two curves, $N(s) = \lambda(s) - \kappa(s)$, and the binormal vector, $B(s) = T(s) \times N(s)$. Using this frame, the twisting number is the average length of the projection of the derivative of the normal vector onto the binormal vector:

$$Tw(\kappa, \lambda) = \frac{1}{2\pi} \int_{s \in \mathbb{S}^1} \langle \dot{N}(s), B(s) \rangle ds.$$

This number may be interpreted as the number of local crossings between κ and λ , counted with a sign and averaged over all directions $u \in \mathbb{S}^2$. To make sense of the idea of a *local crossing*, we use a limit process in which the distance between κ and λ goes to zero. Details are omitted. Similarly, the writhing number of κ may be interpreted as the number of global crossings between κ and λ , again counted with a sign, averaged over all directions, and in the limit when the separation between the knots goes to zero. Since the linking number counts all crossings, we get a relationship between the three measures, which we state without formal proof.

CĂLUGĂREANU-WHITE FORMULA. Let κ be a smooth closed curve in \mathbb{R}^3 and λ one of the two boundary curves of a ribbon centered along κ . Then $Lk(\kappa, \lambda) = Tw(\kappa, \lambda) + Wr(\kappa)$.

Relation to winding number. The writhing number of a space curve is related to the winding number of the curve of critical directions. It is defined such that the directional writhing number remains unchanged as long as we move u on the sphere of directions without crossing the curve and it changes as soon as we cross the curve. The only Reidemeister move that affects the directional writhing number is Type I. The curve of critical directions is therefore traced out by the unit tangent vector and its negative, $T, -T : \mathbb{S}^1 \rightarrow \mathbb{S}^2$. In other words, we have two curves decomposing the sphere into maximal faces of invariant directional writhing number. It will be convenient to identify antipodal points on the sphere and think of a direction as a pair $(u, -u)$. Formally, this means we replace the sphere by the 2-dimensional projective plane, but we don't have to be this formal yet. The pair $(u, -u)$ crosses the curve T iff u crosses T or $-T$.

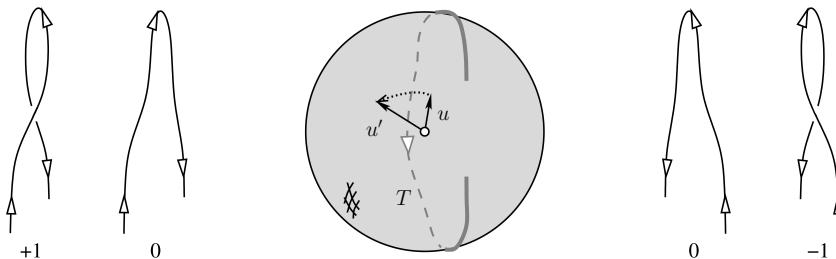


Figure I.11: The change of the viewpoint from u to u' is indicated on the sphere of directions. On the left, this removes a positive crossing, and on the right, this adds a negative crossing. The effect is the same, namely a decrease in the directional writhing number by one. It remains the same even if the curves change their orientation.

Recall that the winding number is defined for a closed curve and a point in the plane. Here we have a closed curve and an antipodal point pair on the sphere. Assuming u and $-u$ are not on the curve, we let the *winding number* be the net number of counterclockwise turns formed by T around the directed line defined by u . We use the same notation as in the plane, denoting this number by $W(T, u)$.

Here we define counterclockwise as seen by looking in the direction u . Figure I.11 illustrates the situation in which $-u$ crosses T from its left to its right. The winding number of T and $(u, -u)$ thus decreases by 1, as the directional writhing number. Indeed, the two change in synchrony in all cases, and we have $DWr(\kappa, u_0) - DWr(\kappa, u) = W(T, u_0) - W(T, u)$ for all $u_0, u \in \mathbb{S}^2$. As a consequence, the average winding number differs from the average directional writhing number by an integer. Integrating the above relation over all directions of the sphere gives $DWr(\kappa, u_0)$ minus $Wr(\kappa)$ on the left and $W(\kappa, u_0)$ minus the average winding number on the right. Hence,

$$Wr(\kappa) = DWr(\kappa, u_0) - W(\kappa, u_0) + \frac{1}{4\pi} \int_{u \in \mathbb{S}^2} W(\kappa, u) du.$$

Bibliographic notes. Knots and links have been studied for centuries, and there are a number of excellent books on the subject, including the text by Adams [2]. Motivation for studying the writhing number of a space curve and the twisting number of a ribbon is derived from the double-helix structure of DNA whose discovery is comparably recent [154]. These numbers measure how wound up, locally and globally, DNA is within the cell [15]. The noteworthy relation between writhing, twisting, and linking numbers has been discovered independently by Călugăreanu [25], Fuller [74], Pohl [119], and White [156]. The relationship to the winding number has been described in [4] and has been used to give an algorithm that computes the writhing number of a closed space polygon in subquadratic time.

I.4 Planar Graphs

Only graphs with relatively few edges can be drawn without crossings in the plane. We consider properties that distinguish such graphs from others. We also prove Tutte's Theorem which implies that every graph that can be drawn without crossing can also be drawn this way with straight edges.

Embeddings. Let $G = (V, E)$ be a simple, undirected graph. A *drawing* maps every vertex $u \in V$ to a point $f(u)$ in \mathbb{R}^2 , and it maps every edge $uv \in E$ to a path with endpoints $f(u)$ and $f(v)$. The drawing is an *embedding* if the points are distinct, the paths are simple and do not cross each other, and incidences are limited to endpoints. It is a *straight-line embedding* if, in addition, all edges are straight line segments. Not every graph can be drawn without crossings. The graph is *planar* if it has an embedding in the plane. As illustrated in Figure I.12 for the complete graph of four vertices, there are many drawings of a planar graph, some with and some without crossings. A *face* of an embedding is a component in the defined decomposition of the plane. We write $n = \text{card } V$, $m = \text{card } E$, and ℓ for the number of faces. Euler's formula is a linear relation between these numbers.

EULER RELATION FOR PLANAR GRAPHS. Every embedding of a connected graph in the plane satisfies $n - m + \ell = 2$.

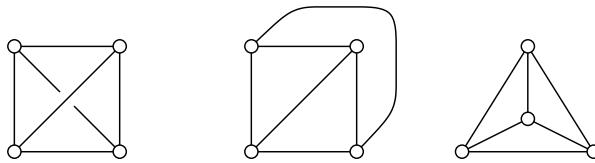


Figure I.12: Three drawings of K_4 . From left to right: a drawing that is not an embedding, an embedding with one curved edge, and a straight-line embedding.

PROOF. Choose a spanning tree of $G = (V, E)$. It has n vertices, $n - 1$ edges, and one face. We have $n - (n - 1) + 1 = 2$, which proves the formula if G is a tree. Otherwise, draw the remaining edges, one at a time. Each edge decomposes one face into two, thus maintaining the relation by increasing both the number of edges and the number of faces by one. \square

If the graph has more than one connected component, then the right-hand side of the equation is replaced by one plus that number. Note that the Euler Relation implies that the number of faces is the same for all embeddings and is therefore a property of the graph. We get bounds on the number of edges and faces, in terms of the number of vertices, by considering *maximally connected* planar graphs for which adding any one edge would violate planarity. Every face of a maximally connected planar graph with three or more vertices is necessarily a triangle, for if there is a face with more than three edges, we can add a path that crosses none of the earlier paths. Let $n \geq 3$ be the number of vertices, as before. Since every face has three edges and every edge belongs to two triangles, we have $3\ell = 2m$. We use this relation to rewrite the Euler Relation: $n - m + \frac{2m}{3} = 2$ and $n - \frac{3\ell}{2} + \ell = 2$ and hence $m = 3n - 6$ and $\ell = 2n - 4$. Every planar graph can be completed to a maximally connected planar graph, which implies that it has at most these numbers of edges and faces.

Non-planarity. We can use the Euler Relation to prove that the complete graph of five vertices and the complete bipartite graph of three plus three vertices are not planar. Consider first K_5 , which is drawn in Figure I.13, left. It has $n = 5$ vertices and $m = 10$ edges, contradicting the upper bound of at most $3n - 6 = 9$ edges

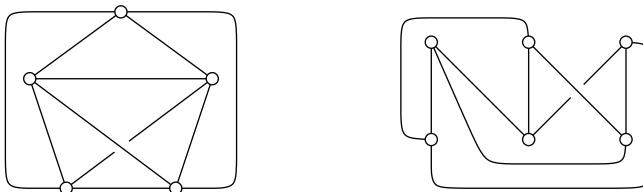


Figure I.13: K_5 on the left and $K_{3,3}$ on the right, each drawn with the unavoidable one crossing.

for maximally connected planar graphs. Consider second $K_{3,3}$, which is drawn in Figure I.13, right. It has $n = 6$ vertices and $m = 9$ edges. Each cycle has even length, which implies that each face of a hypothetical embedding has four or more edges. We get $4\ell \leq 2m$ and $m \leq 2n - 4 = 8$ after plugging the inequality into the Euler Relation, again a contradiction.

In a sense, K_5 and $K_{3,3}$ are the quintessential non-planar graphs. Two graphs are *homeomorphic* if one can be obtained from the other by a sequence of operations, each deleting a degree-2 vertex and merging their two edges into one or doing the inverse.

KURATOWSKI THEOREM. A simple graph is planar iff no subgraph is homeomorphic to K_5 or to $K_{3,3}$.

The proof of this result is omitted. The remainder of this section focuses on straight-line embeddings of planar graphs.

Convex combinations. Two points $a_0 \neq a_1$ define a unique line that passes through both. Each point on this line can be written as $x = (1-t)a_0 + ta_1$, for some $t \in \mathbb{R}$. For $t = 0$ we get $x = a_0$, for $t = 1$ we get $x = a_1$, and for $0 < t < 1$ we get a point in between. If we have more than two points, we repeat the construction by adding all points $y = (1-t)x + ta_2$ for which $0 \leq t \leq 1$, and so on, as illustrated in Figure I.14. Given $k+1$ points a_0, a_1, \dots, a_k , we can do the same construction in

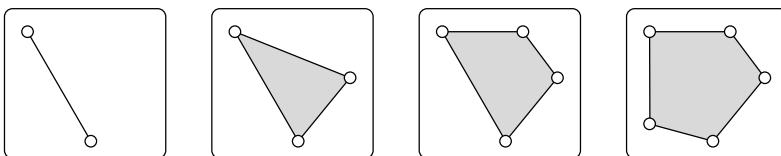


Figure I.14: From left to right: the construction of the convex hull of five points by adding one point at a time.

one step, calling a point $x = \sum_{i=0}^k t_i a_i$ a *convex combination* of the a_i if $\sum_{i=0}^k t_i = 1$ and $t_i \geq 0$ for all $0 \leq i \leq k$. The set of convex combinations is the *convex hull* of the a_i .

We are interested in graphs that arise as edge-skeletons of triangulations of the disk, like the one in Figure I.15. Letting $G = (V, E)$ be such a graph, we distinguish edges and vertices on the boundary from the ones in the interior of the disk. When we embed G in \mathbb{R}^2 , we make sure that the boundary edges and vertices map to the boundary of the outer face. Since we only consider straight-line embeddings, it suffices to study mappings of the vertex set into the plane. We call $f : V \rightarrow \mathbb{R}^2$ a *strictly convex combination mapping* if for every interior vertex $u \in V$ there are real numbers $t_{uv} > 0$ with $\sum_v t_{uv} = 1$ and $f(u) = \sum_v t_{uv} f(v)$, where both sums are over all neighbors v of u . In words, every interior vertex maps to a point in the interior of the convex hull of the images of its neighbors. We will repeatedly use this mapping

in combination with a linear function $h : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by $h(x) = \langle x, p \rangle + c$, where $p \in \mathbb{R}^2$ is a non-zero vector and c is a real number. Composing f with h , we get $h(f(u)) = \sum_v t_{uv}h(f(v))$. In words, the value we get for u is the same strictly convex combination of the values for its neighbors.

Straight-line embeddings. Suppose we have a straight-line embedding of G in which the boundary edges map to the boundary of the outer face. Then every interior vertex lies inside the cycle connecting its neighbors. It follows that this embedding defines a strictly convex combination mapping. We now show that the reverse is also true provided the boundary vertices map to the corners of a strictly convex polygon.

TUTTE'S THEOREM. Let $G = (V, E)$ be the edge-skeleton of a triangulation of the disk and $f : V \rightarrow \mathbb{R}^2$ a strictly convex combination mapping that maps the boundary vertices to the corners of a strictly convex polygon. Then drawing straight edges between the images of the vertices gives a straight-line embedding.

We will give the proof in three steps, which we now prepare with two observations. A *separating edge* of G is an interior edge that connects two boundary vertices. It is convenient to assume that G has no separating edge, but if it does, we can split the graph into two and do the argument for each piece. Call a path in G *interior* if all its vertices are interior except possibly the first and the last. Under the assumption of no separating edge, every interior vertex u can be connected to every boundary vertex by an interior path. Indeed, we can find an interior path that connects u to a first boundary vertex w . Let w_0 and w_1 be the neighboring boundary vertices. Since none of the edges separate, the neighbors of w form a unique interior path connecting w_0 to w_1 . It follows that there is an interior path connecting u to w_0 . By repeating the argument, substituting w_0 for w , we eventually see that u has interior paths to all boundary vertices.

Secondly, suppose that $h \circ f$ takes its maximum at an interior vertex, u . Since $h \circ f(u)$ is a strictly convex combination of the values the neighbors, we conclude that $h \circ f(v) = h \circ f(u)$ for all neighbors v of u . We can iterate and because of the mentioned interior path property we eventually reach every vertex. It thus follows that $h \circ f$ has the same value at all vertices of G . We refer to this observation as the *maximum principle* and to its symmetric version as the *minimum principle*.

Proof of Tutte's Theorem. We now present the proof in three steps. First, all interior vertices u of V map to the interior of the strictly convex polygon whose corners are the images of the boundary vertices. To see this, choose $p \in \mathbb{R}^2$ and $c \in \mathbb{R}$ such that the line $h^{-1}(0)$ defined by $h(x) = \langle x, p \rangle + c$ contains a boundary edge and $h(f(w)) > 0$ for all boundary vertices other than the endpoints of that edge. Then $h(f(u)) > 0$; otherwise, the minimum principle would imply $h(f(v)) = 0$ for all vertices. Repeating this argument for all edges of the convex polygon implies that all interior vertices u have $f(u)$ in the interior of the polygon. This implies in

particular that each triangle incident to a boundary edge is non-degenerate, that is, its three vertices are not collinear.

Second, letting yuv and zuv be the two triangles sharing the interior edge uv in G , the points $f(y)$ and $f(z)$ lie on opposite sides of the line $h^{-1}(0)$ that passes through $f(u)$ and $f(v)$. To see this, assume $h(f(y)) > 0$ and find a strictly rising path connecting y to the boundary. It exists because $h(f(y)) > h(f(u))$, so one of the neighbors of y has strictly larger function value, and the same is true for the next vertex on the path, and so on. Since $h(f(y)) > 0$, both u and v have neighbors for which $h \circ f$ is negative. So we can find a strictly falling path connecting u to the boundary and the same for v , as illustrated in Figure I.15. The rising path

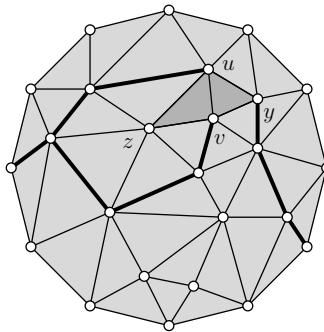


Figure I.15: One strictly rising and two strictly falling paths connecting y , u , and v to the boundary.

does not cross the falling paths, but the two falling paths may share a vertex, as in Figure I.15. In either case, we get a piece of the triangulation bounded by vertices with non-positive function values. Other than u and v , all other vertices in this boundary have strictly negative function values. If z belongs to the boundary of this piece, then it has strictly negative function value simply because it differs from u and v . Otherwise, it belongs to the interior of the piece, and we have $h(f(z)) < 0$ by the maximum principle. We note that this argument uses $h(f(y)) > 0$ in an essential manner. To show that this assumption is justified, we connect yuv by a sequence of triangles until we reach a boundary edge. In this sequence, any two contiguous triangles share an edge. As observed in the first step, the image of the last triangle is non-degenerate. Going backward, this implies that the image of the second to the last triangle is non-degenerate, and so on. Finally, the image of yuv is non-degenerate, as required.

Third, no two of the edges cross. To get a contradiction, assume x is a point in the common interiors of two edges, uv and $u'v'$. Choose a half-line that emanates from x and avoids the images of all vertices. Since the vertices y and z that form triangles with uv map to opposite sides of the line passing through $f(u)$ and $f(v)$, the half-line intersects exactly one of the edges yu , yv , zu , zv . Continuing along the half-line, we get a sequence of edges starting with uv and ending with a boundary edge. Similarly, the half-line defines another sequence of edges starting with $u'v'$

and ending with the same boundary edge. Going back in both sequences, we pass from one edge to an unambiguously defined preceding edge. Since we start with the same boundary edge, we get $uv = u'v'$. This completes the proof of Tutte's Theorem.

Constructing straight-line embeddings. Tutte's Theorem leads to a simple algorithm for constructing a straight-line embedding of a planar graph. For simplicity, we assume that it is the edge-skeleton of a triangulation of the disk and that none of its edges separates. We reindex such that u_1 to u_k are ordered along the boundary of the outer face and u_{k+1} to u_n are the interior vertices of the graph. First, we set $f(u_i) = (\cos(2i\pi/k), \sin(2i\pi/k))$, for $1 \leq i \leq k$, to place the boundary vertices in order on the unit circle in the plane. They form the corners of a strictly convex polygon, as required. Expressing the image of each interior vertex as a strictly convex combination of the images of its neighbors, we write

$$f(u_j) = \frac{1}{d_j} \sum f(v),$$

for each $k+1 \leq j \leq n$, where d_j is the degree of u_j and the sum is over all neighbors v of u_j in the graph. We get a system of $n - k$ linear equations in $n - k$ unknowns, the images of the interior vertices. Writing the system in matrix form, we get one non-zero coefficient for each interior vertex and two more for each edge connecting two interior vertices. By the Euler Relation, the number of edges is less than $3n$. It follows that the system is sparse with fewer than $7n$ non-zero coefficients. It thus permits efficient methods to find the solution, which by Tutte's Theorem corresponds to a straight-line embedding of the graph.

Bibliographic notes. Graphs that can be drawn in the plane without crossings arise in a number of applications, including geometric modeling, geographic information systems, and others. We refer to [117] for a collection of mathematical and algorithmic results specific to planar graphs. The fact that all planar graphs have straight-line embeddings has been known long before Tutte's Theorem. Early last century, Steinitz showed that every 3-connected planar graph is the edge-skeleton of a convex polytope in \mathbb{R}^3 [137]. This skeleton can be projected to \mathbb{R}^2 to give a straight-line embedding. In the 1930s, Koebe proved that every planar graph is the intersection graph of a collection of possibly touching but not otherwise overlapping closed disks in \mathbb{R}^2 [95]. We get a straight-line embedding by connecting the centers of the touching disks. The original theorem by Tutte is for coefficients t_{uv} equal to one over the degree of u [141]. The more general version and the proof presented in this section are fashioned after the more recent paper by Floater [68]. Efficient numerical methods for solving systems of linear equations can be found in the linear algebra text by Strang [139].

Exercises

The credit assignment reflects a subjective assessment of difficulty. A typical question can be answered using knowledge of the material combined with some thought and analysis.

1. **Deciding connectivity** (two credits). Given a simple graph with n vertices and m edges, the disjoint set system takes time proportional to $(n + m)\alpha(n)$ to decide whether or not the graph is connected.
 - (i) Describe a different algorithm that makes the same decision in time proportional to $n + m$.
 - (ii) Modify the algorithm so it computes the connected components in time proportional to $n + m$.
2. **Shelling disks** (two credits). Consider a triangulation of a simple closed polygon in the plane, but one that may have interior vertices inside the polygon. A *shelling* is a total order of the triangles such that the union of the triangles in any initial sequence is homeomorphic to a closed disk. Prove that every such triangulation has a shelling.
3. **Jordan curve** (one credit). Recall the Jordan Curve Theorem, which says that every simple closed curve in the plane decomposes \mathbb{R}^2 into two components.
 - (i) Show the same is true for a simple closed curve on the sphere, $\mathbb{S}^2 = \{x \in \mathbb{R}^3 \mid \|x\| = 1\}$.
 - (ii) Give an example that shows the result does not hold for simple closed curves on the torus.
4. **Homeomorphisms** (one credit). Give explicit homeomorphisms to show that the following spaces with topologies inherited from the respective Euclidean spaces that contain them are homeomorphic:
 - $\mathbb{R}^1 = \mathbb{R}$, the real line;
 - $(0, 1)$, the open interval;
 - $\mathbb{S}^1 - \{(0, 1)\}$, the circle with one point removed.
 Generalize your homeomorphisms to show the same for the Euclidean plane, the open disk, and the sphere with one point removed.
5. **Splitting a link** (two credits). Prove that the Borromean rings shown in Figure I.16 on the left are not splittable.
6. **Deforming a link** (two credits). Use Reidemeister moves to demonstrate that the two links in Figure I.16 in the middle and on the right are equivalent.
7. **Planar graph coloring** (two credits). Recall that every planar graph has a vertex of degree at most five. We can use this fact to show that every planar graph has a vertex 6-coloring, that is, a coloring of each vertex with one of six colors such that any two adjacent vertices have different colors. Indeed, after removing a vertex with fewer than six neighbors we use induction to 6-color

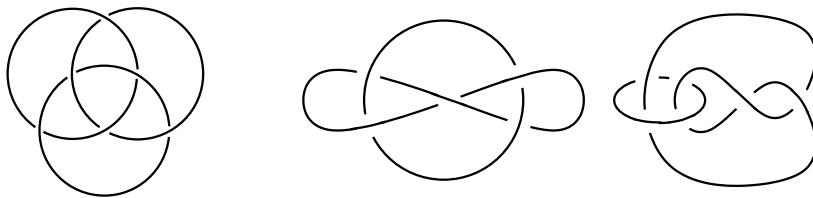


Figure I.16: Left: any two of the three knots of the Borromean rings can be split but are held together by the third knot. Right: two generic projections of the Whitehead link.

the remaining graph, and when we put the vertex back, we choose a color that differs from the colors of its neighbors. Refine the argument to prove that every planar graph has a vertex 5-coloring.

8. **Edge coloring** (three credits). We color each edge of a maximally connected planar graph with one of three colors such that each face (triangle) has all three colors in its boundary.
 - (i) Show that a 4-coloring of the vertices implies a 3-coloring of the edges.
 - (ii) Show that a 3-coloring of the edges implies a 4-coloring of the vertices.

In other words, proving that every planar graph has a vertex 4-coloring is equivalent to proving that every triangulation in the plane has an edge 3-coloring.

Chapter II

Surfaces

The most common 2-dimensional spaces are 2-manifolds, or surfaces, which come in two varieties: with and without boundary. We usually envision them put into 3-dimensional space, sometimes with and preferably without self-intersections. Not all surfaces can be embedded in 3-dimensional Euclidean space and self-intersections are unavoidable, but often they are accidental. Indeed, choosing a nice embedding of a surface in space is an interesting computational problem. We address this question for surfaces made out of triangles.

II.1 2-dimensional Manifolds

In our physical world, the use of the term surface usually implies a 3-dimensional, solid shape of which this surface is the boundary. In mathematics, the solid shape is not assumed, and we discuss surfaces in their own right. Indeed, there are closed surfaces that are not the boundary of any solid shape. They are non-orientable and do not embed into 3-dimensional Euclidean space, which is why our intuition for them is lacking.

Topological 2-manifolds. Consider the open disk of points at distance less than one from the origin, $D = \{x \in \mathbb{R}^2 \mid \|x\| < 1\}$. It is homeomorphic to \mathbb{R}^2 , as for example established by the homeomorphism $f : D \rightarrow \mathbb{R}^2$ defined by $f(x) = x/(1 - \|x\|)$. We will call any subset of a topological space that is homeomorphic to D an open disk. A *2-manifold (without boundary)* is a topological space M whose points all lie in open disks. Intuitively, this means that M looks locally like the plane.

M is *compact* if for every covering of M by open sets, called an *open cover*, we can find a finite number of the sets that cover M . We say that the open cover always has *finite subcover*. Examples of non-compact 2-manifolds are \mathbb{R}^2 itself and open subsets of \mathbb{R}^2 . Examples of compact 2-manifolds are shown in Figure

II.1, top row. We get *2-manifolds with boundary* by removing open disks from 2-manifolds without boundary. Alternatively, we could require that each point have a neighborhood homeomorphic to either D or D_+ , the half-disk obtained by removing all points with negative second coordinate from D . The *boundary* of a 2-manifold with boundary consists of all points x whose neighborhoods are homeomorphic to D_+ . Within the boundary, the neighborhood of every point x is an open interval, which is the defining property of a *1-manifold*, or *curve*. There is only one type of connected, compact 1-manifold, namely the closed curve. Following the practice of considering topologically equivalent spaces the same, we will therefore often refer to a closed curve as a circle. If M is compact, this implies that its boundary is a collection of circles. Examples of 2-manifolds with boundary are the (closed) disk, the cylinder, and the Möbius strip, all illustrated in Figure II.1, bottom row.

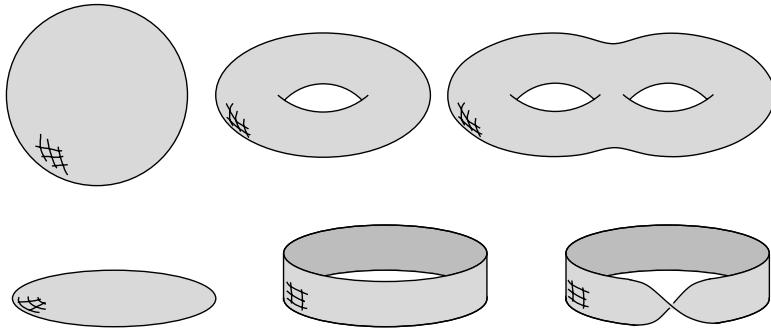


Figure II.1: Top from left to right: the sphere, S^2 , the torus, \mathbb{T}^2 , the double torus, $\mathbb{T}^2 \# \mathbb{T}^2$. Bottom from left to right: the disk, the cylinder, the Möbius strip.

We get new 2-manifolds from old ones by gluing them to each other. Specifically, remove an open disk from each of two 2-manifolds, M and N , find a homeomorphism between the two boundary circles, and identify corresponding points. The result is the *connected sum* of the two manifolds, denoted as $M \# N$. Forming the connected sum with the sphere does not change the manifold since it just means replacing one disk by another. Adding the torus is the same as attaching the cylinder at both boundary circles after removing two open disks.

Orientability. Of the examples we have seen so far, the Möbius strip has the curious property that it seems to have two sides locally at every interior point but there is only one side globally. To express this property intrinsically, that is, without reference to the embedding in \mathbb{R}^3 , we consider a small, oriented circle inside the strip. We move it around without altering its orientation, like a clock whose fingers keep turning in the same direction. However, if we slide the clock once around the strip, its orientation is the reverse of what it used to be, and we call the path of its center an *orientation-reversing* closed curve. There are also *orientation-preserving* closed curves in the Möbius strip, such as the one that goes around the

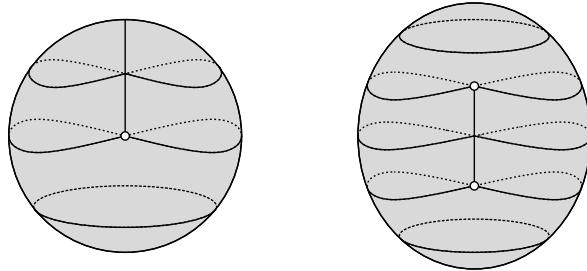


Figure II.2: Left: the projective plane, \mathbb{P}^2 , obtained by gluing a disk to a Möbius strip. Right: the Klein bottle obtained by gluing two Möbius strips together. The vertical lines are self-intersections that we ignore.

strip twice following along close to the boundary. If all closed curves in a 2-manifold are orientation-preserving, then the 2-manifold is *orientable*; otherwise, it is *non-orientable*. The curves drawn on the projective plane and the Klein bottle in Figure II.2 are all orientation-preserving. We leave finding orientation-reversing curves on the same two surfaces as an instructive exercise for the reader.

Note that the boundary of the Möbius strip is a single circle. We can therefore glue the strip to a sphere or a torus after removing an open disk from the latter. This operation is often referred to as adding a *cross-cap* to the sphere or torus. The result is homeomorphic to the surface obtained by identifying antipodal points on the circle where the disk was removed. In the case of the sphere, we get the *projective plane*, the sphere with one cross-cap, and in the case of the torus, we get the *Klein bottle*, the sphere with two cross-caps. Both cannot be embedded in \mathbb{R}^3 , so we have to draw them with self-intersections, but these should be ignored when we think about these surfaces.

Classification. As it turns out, we have seen examples of each major kind of compact 2-manifold. They were completely classified about a century ago by cutting and gluing to arrive at a unique representation for each type. This representation is a convex polygon whose edges are glued in pairs, called a *polygonal schema*. Figure II.3 shows that the sphere, the torus, the projective plane, and the Klein bottle can all be constructed from the square. More generally, we have a $4g$ -gon for a sphere with g tubes and a $2g$ -gon for a sphere with g cross-caps attached to it. The gluing pattern is shown in the second row of Figure II.3. Note that the square of the torus is in standard form but that of the Klein bottle is not.

CLASSIFICATION THEOREM FOR COMPACT 2-MANIFOLDS. The two infinite families $\mathbb{S}^2, \mathbb{T}^2, \mathbb{T}^2 \# \mathbb{T}^2, \dots$ and $\mathbb{P}^2, \mathbb{P}^2 \# \mathbb{P}^2, \dots$ exhaust the compact 2-manifolds without boundary.

The first family of orientable, compact 2-manifolds consists of the sphere, the torus, the double torus, and so on. The second family of non-orientable, compact 2-manifolds consists of the projective plane, the Klein bottle, the triple projective

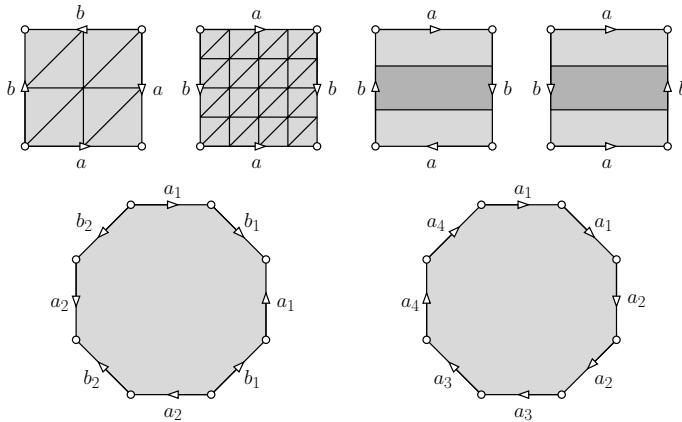


Figure II.3: Top from left to right: the sphere, the torus, the projective plane, and the Klein bottle. After removing the (darker) Möbius strip from the last two, we are left with a disk in the case of the projective plane and another Möbius strip in the case of the Klein bottle. Bottom: the polygonal schema in standard form for the double torus on the left and the double Klein bottle on the right.

plane, and so on. To get a classification of the connected, compact 2-manifolds with boundary, we can take one without boundary and make h holes by removing the same number of open disks. Each starting compact 2-manifold and each $h \geq 1$ give a different surface, and they exhaust all possibilities.

Triangulations. To triangulate a 2-manifold, we decompose it into triangular regions, each a disk whose boundary circle is cut at three points into three paths. We may think of the region and its boundary as the homeomorphic image of a triangle. By taking a geometric triangle for each region and arranging them so they share vertices and edges the same way as the regions, we obtain a piecewise linear model which is a *triangulation* if it is homeomorphic to the 2-manifold. See Figure II.4 for a triangulation of the sphere. Since the triangles are geometric, the condition of homeomorphism requires that any two either be disjoint, share an edge, or share a vertex. Sharing two edges is not permitted for then the two triangles would be the same. It is also not permitted that two vertices of a triangle be the same. To illustrate these conditions, we note that the triangulation of the first square in Figure II.3 is not a valid triangulation of the sphere, but the triangulation of the second square is a valid triangulation of the torus.

Given a triangulation of a 2-manifold \mathbb{M} , we may orient each triangle. Two triangles sharing an edge are *consistently oriented* if they induce opposite orientations on the shared edge, as in Figure II.4. Then \mathbb{M} is orientable iff the triangles can be oriented in such a way that every adjacent pair is consistently oriented.

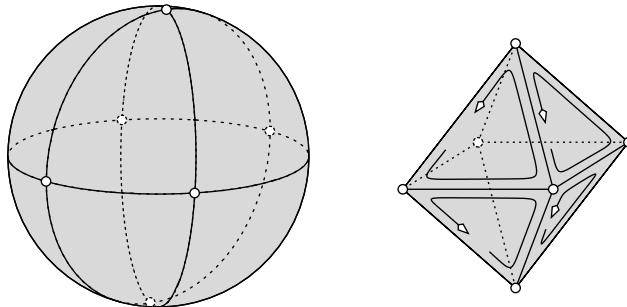


Figure II.4: The sphere is homeomorphic to the surface of an octahedron, which is a triangulation of the sphere.

Euler characteristic. Recall that a triangulation is a collection of triangles, edges, and vertices. We are only interested in finite triangulations. Letting n , m , and ℓ be the numbers of vertices, edges, and triangles, as in the previous chapter, the *Euler characteristic* is their alternating sum, $\chi = n - m + \ell$. We have seen that the Euler characteristic of the sphere is $\chi = 2$, no matter how we triangulate. More generally, the Euler characteristic is independent of the triangulation for every 2-manifold.

EULER CHARACTERISTIC OF COMPACT 2-MANIFOLDS. A sphere with g tubes has $\chi = 2 - 2g$ and a sphere with g cross-caps has $\chi = 2 - g$.

The number g is the *genus* of \mathbb{M} ; it is the maximum number of disjoint closed curves along which we can cut without disconnecting \mathbb{M} . To see this result, we may triangulate the polygonal schema of \mathbb{M} . For a sphere with g tubes we have $\ell = 1$ region, $m = 2g$ edges, and $n = 1$ vertex. Further decomposing the edges and regions does not change the alternating sum, so we have $\chi = 2 - 2g$. For a sphere with g cross-caps we have $\ell = 1$ region, $m = g$ edges, and $n = 1$ vertex giving $\chi = 2 - g$.

Observe that adding a tube decreases the Euler characteristic by two, while adding a cross-cap decreases it by only one. Indeed, we can substitute k handles for $2k$ cross-caps and obtain the g -fold projective plane from the k -fold torus by gluing $g - 2k$ cross-caps, provided $g > 2k$. Note that non-orientability cannot be cancelled by the connected sum. Hence, this operation can get us from the orientable to the non-orientable manifolds but not back.

Doubling. The compact, non-orientable 2-manifolds can be obtained from the orientable 2-manifolds by identifying points in pairs. For example, if we identify opposite (*antipodal*) points of the sphere, we get the projective plane. We can also go in the other direction, constructing orientable manifolds from non-orientable ones; see Figure II.5. Imagine a triangulation of a connected, compact, non-orientable 2-manifold \mathbb{N} in \mathbb{R}^3 , drawn with self-intersections, which we ignore. Make two copies of each triangle, edge, and vertex offsetting them slightly, one on either side of the manifold. Here sidedness is local and therefore well defined. The triangles

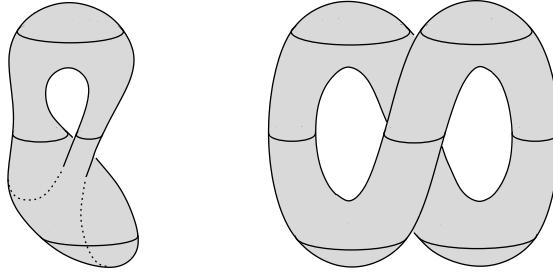


Figure II.5: Doubling the Klein bottle produces the torus.

fit together locally, and because N is connected and non-orientable, they form the triangulation of a connected 2-manifold, M . It is orientable because one side is consistently facing N . Since all triangles, edges, and vertices are doubled, we have $\chi(M) = 2\chi(N)$. Using the relation between genus and Euler characteristic, we have $\chi(N) = 2 - g(N)$ and therefore $\chi(M) = 4 - 2g(N) = 2 - 2g(M)$. It follows that M has $g(M) = g(N) - 1$ tubes. As listed in Table II.1, the doubling operation constructs the sphere from the projective plane, the torus from the Klein bottle, etc. The result of the doubling operation is sometimes called the *double cover*, since the reverse operation of re-identifying doubled regions maps M to N , covering it twice.

$\chi(N)$	$g(N)$	N	M	$g(M)$	$\chi(M)$
1	1	\mathbb{P}^2	S^2	0	2
0	2	$\mathbb{P}^2 \# \mathbb{P}^2$	T^2	1	0
-1	3	$\mathbb{P}^2 \# \mathbb{P}^2 \# \mathbb{P}$	$T^2 \# T^2$	2	-2
...

Table II.1: Doubling turns the non-orientable 2-manifold on the left into the orientable 2-manifold on the right.

Bibliographic notes. The confusing aspects of non-orientable 2-manifolds have been captured in a delightful novel about the life within such a surface [1]. The classification of compact 2-manifolds is sometimes credited to Brahana [20] and at other times to Dehn and Heegard [43]. The classification of 3-manifolds, on the other hand, is an ongoing project within mathematics. With the proof of the Poincaré conjecture by Perelman, there is new hope that the classification of 3-manifolds can be accomplished soon. In contrast, recognizing whether two triangulated 4-manifolds are homeomorphic is undecidable [104]. The classification of manifolds from first principles beyond dimension three is therefore a hopeless undertaking.

II.2 Searching a Triangulation

Many algorithms benefit from a convenient data structure that represents a surface by storing its triangulation. In this section, we describe such a data structure and show how to use it to determine the topological type of a surface.

Ordered triangles. We begin with the description of the core piece of the data structure, which is a representation of the symmetry group of the standard triangle. Its main function will be to keep track of direction and orientation when we navigate the triangulation. This group is isomorphic to the group of permutations of three elements, the vertices of the triangle. We call each permutation an *ordered triangle* and use cyclic shifts and transpositions to move between them. As illustrated in

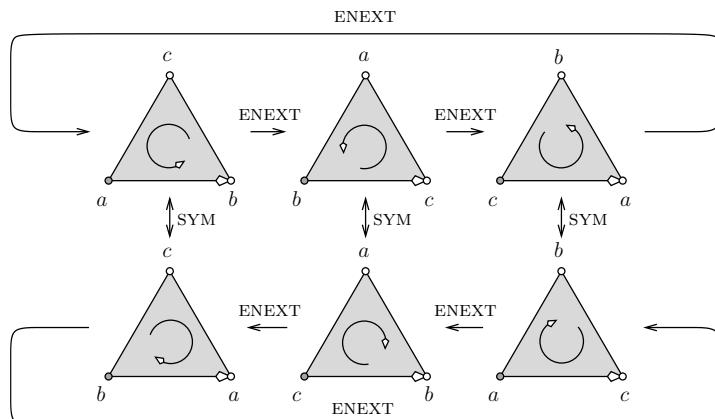


Figure II.6: The symmetry group of the standard triangle consists of six ordered versions. The cyclic shifts partition the group into two orientations, each consisting of three ordered triangles.

Figure II.6, the cyclic shift from abc to bca corresponds to advancing the leading directed edge to the next position, from ab to bc . The transposition of the leading two vertices corresponds to reversing the direction of the lead edge while keeping the third vertex fixed.

We store each triangle in a single node of the data structure, to be described shortly. A reference to the triangle consists of a pointer to this node, μ , together with a three-bit integer, ι , identifying the ordered version of the triangle. Using the first bit to identify the orientation, we represent $abc, bca, cab, bac, cba, acb$ by $\iota = 0, 1, 2, 4, 5, 6$, in this sequence. Moving between different ordered versions of the same triangle can be done with simple arithmetic operations on ι . To advance the lead edge, we increment using modulo arithmetic.

```

ordTri ENEXT( $\mu, \iota$ )
  if  $\iota \leq 2$  then return  $(\mu, (\iota + 1) \bmod 3)$ 
  else return  $(\mu, (\iota + 1) \bmod 3 + 4)$ 
endif.

```

To reverse the direction of the lead edge, we flip the first bit.

```

ordTri SYM( $\mu, \iota$ )
  return  $(\mu, (\iota + 4) \bmod 8)$ .

```

We see that encoding the symmetry group requires very little overhead, just a few bits whenever we point to a triangle.

Data structure. We are now ready to describe the data structure representing the triangulation K of a connected, compact, 2-manifold without boundary. We store the vertices of K in a linear array, $V[1..n]$. We store the triangles in the nodes of a graph, by which we mean a data structure consisting of memory locations with pointers referring to each other. The arcs connect nodes of neighboring triangles defined by shared edges. Since every triangle has exactly three neighbors, the degree of every node is three. Inside a node, we store pointers to the three neighbors as well as to the three vertices, which are indices into V .

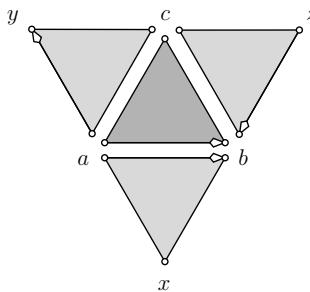


Figure II.7: The triangle abc with its three neighbors. The arrowheads identify the directed lead edges.

Let abc be a triangle and x, y, z the respective third vertices of the neighbor triangles. Each ordered version of the triangle points to its lead vertex and the ordered neighbor triangle that shares the directed lead edge. To describe this in an example, we assume the nodes μ, μ_x, μ_y, μ_z store the four triangles with $\iota = 0$ corresponding to the ordered versions abc, abx, ayc, zbc , as drawn in Figure II.7. Assuming a is stored at positions i in V and observing that ab is the lead edge of abx , the ordered triangle abc stores pointers $(\mu, 0).org = i$ and $(\mu, 0).fnext = (\mu_x, 0)$. Assuming furthermore that b and c are stored at positions j and k of the vertex array, the other five ordered triangles in μ store pointers to the positions j, k, j, k, i and to the ordered triangles $(\mu_z, 1), (\mu_y, 2), (\mu_x, 4), (\mu_z, 5), (\mu_y, 6)$, in this

sequence. To move around in the triangulation, we use simple functions to retrieve this information.

```
ordTri FNEXT( $\mu, \iota$ )
    return ( $\mu, \iota$ ).fnext.
```

```
int ORG( $\mu, \iota$ )
    return ( $\mu, \iota$ ).org.
```

There is clearly redundancy left in the proposed data structure, but we resist further optimizations to keep the implementation transparent.

Depth-first Search. A common operation is visiting all triangles of the triangulation. This corresponds to searching the entire representing graph. Two of the most popular strategies are Breadth-first Search and Depth-first Search. As suggested by the name, Breadth-first Search proceeds along an advancing front that expands around an initial node. In contrast, Depth-first Search ventures directly into the unknown and covers the neighborhood only after returning from the adventure. We implement the latter strategy using a recursive function. Assuming all nodes are initially unmarked, we start the search by calling that function for an arbitrary first node μ_0 .

```
void VISIT( $\mu$ )
    if  $\mu$  is unmarked then mark  $\mu$ ; P1;
        forall neighbors  $\nu$  of  $\mu$  do
            VISIT( $\nu$ )
        endfor; P2
    else P3
    endif.
```

The search proceeds along a spanning tree of the graph defined by calling a neighboring node ν a *child* of μ if the first visit to ν originates from μ . The root of this tree is μ_0 . To customize the function, we would add instructions at the three indicated places:

P1: steps to be executed the first time the node is visited;

P2: steps to be executed after all children have been processed;

P3: steps to be executed each time the node is revisited.

We will see examples of such customizations shortly. After searching the graph once, we will typically search it once more to remove all the marks and prepare the graph for further processing. Without accounting for the additional instructions, the running time of Depth-first Search is linear in $n + m$, the number of nodes and arcs in the graph. Indeed, each arc is traversed exactly twice, once in each direction.

Orientability. We use Depth-first Search to decide whether a connected, compact 2-manifold without boundary given by a triangulation K is orientable. We do this by orienting all triangles in a consistent manner and report non-orientability if the attempt fails. In other words, we choose one of two orientations for each triangle such that the shared edges between neighboring triangles are directed in opposite ways. Assuming none of the orientations are yet chosen, we start the process by calling the function for an arbitrary first ordered triangle (μ_0, ι_0) .

```

boolean ISORIENTABLE( $\mu, \iota$ )
  if  $\mu$  is unmarked then mark  $\mu$  and choose orientation containing  $\iota$ ;
     $b_x = \text{ISORIENTABLE}(\text{FNEXT}(\text{SYM}(\mu, \iota)))$ ;
     $b_y = \text{ISORIENTABLE}(\text{FNEXT}(\text{ENEXT}(\text{SYM}(\mu, \iota))))$ ;
     $b_z = \text{ISORIENTABLE}(\text{FNEXT}(\text{ENEXT}^2(\text{SYM}(\mu, \iota))))$ ;
    return  $b_x$  and  $b_y$  and  $b_z$ 
  else return [orientation of  $\mu$  contains  $\iota$ ]
  endif.

```

Here we orient μ at P1, we unwind the for-loop, and we return a boolean value at P2 and another at P3. The latter value indicates whether or not we have consistent orientations in spite of the triangle μ having been oriented prior to the current visit. The boolean value returned at P2 indicates whether or not we have found a contradiction to orientability. A single value of FALSE anywhere during the computation is propagated to the root of the search tree, telling us that the surface is non-orientable. Since each triangle has only three neighbors, the running time of the algorithm is linear in the number of triangles.

Classification. Recall from the preceding section that the type of a connected, compact 2-manifold without boundary is uniquely determined by its genus and whether or not it is orientable. Since every triangle has three edges and every edge belongs to two triangles, we have $3\ell = 2m$ and therefore $2n - \ell = 4 - 4g$ in the orientable case and $2n - \ell = 4 - 2g$ in the non-orientable case. Assuming we know the number of vertices from the size of the array, we just need to count the triangles, which we do again by Depth-first Search.

```

int #TRIANGLES( $\mu, \iota$ )
  if  $\mu$  is unmarked then mark  $\mu$ ;
     $\ell_x = \#TRIANGLES(\text{FNEXT}(\mu, \iota))$ ;
     $\ell_y = \#TRIANGLES(\text{FNEXT}(\text{ENEXT}(\mu, \iota)))$ ;
     $\ell_z = \#TRIANGLES(\text{FNEXT}(\text{ENEXT}^2(\mu, \iota)))$ ;
    return  $\ell_x + \ell_y + \ell_z + 1$ 
  else return 0
  endif.

```

Combining the information, it is now easy to determine the genus.

```

int GENUS( $\mu, \iota$ )
 $\ell = \#$ TRIANGLES( $\mu, \iota$ );
if ISORIENTABLE( $\mu, \iota$ ) then return  $(\ell - 2n + 4)/4$ 
else return  $(\ell - 2n + 4)/2$ 
endif.

```

In summary, we can decide the topological type of a triangulated, compact 2-manifold without boundary in time linear in the number of triangles. Clearly, we cannot do it faster since the entire triangulation must be searched; otherwise, we could alter the type by a small modification. By adding another search counting the boundaries, we can extend this result to compact 2-manifolds with boundary.

Bibliographic notes. Data structures for storing triangulated 2-manifolds have been described in the computer science literature since Baumgart [16]; see also the doubly-linked edge lists in [124] and the quad-edge structure in [80]. These data structures differ in their details from the graph representation described in this section but are functionally very similar. Extensions to storing 3- and higher-dimensional complexes can be found in [51] and in [21]. Searching graphs is a core topic in computer science, and descriptions of Depth-first Search can be found in most algorithms texts, including [6] and [41].

II.3 Self-intersections

Since non-orientable, compact 2-manifolds without boundary cannot be embedded in 3-dimensional Euclidean space, all their models in that space occur with self-intersections. In contrast, all orientable, compact 2-manifolds have embeddings, but their models may have accidental self-intersections. Removing these self-intersections is a core topic in repairing surface models of solid shapes.

Mapping into space. Let \mathbb{M} be a compact 2-manifold without boundary. We want to say what it means for \mathbb{M} to be smooth and for a continuous map $f : \mathbb{M} \rightarrow \mathbb{R}^3$ to be a smooth mapping. We define a *coordinate chart*, (U, ϕ) , to be an open set, $U \subset \mathbb{M}$, together with a continuous map, $\phi : U \rightarrow \mathbb{R}^2$, that is a homeomorphism onto its image. Two coordinate charts, (U, ϕ) and (V, ψ) , are *compatible* if U and V are disjoint, or the map

$$\phi \circ \psi^{-1} : \psi(U \cap V) \rightarrow \phi(U \cap V)$$

extends to a smooth map from \mathbb{R}^2 to \mathbb{R}^2 . The class of such functions is often referred to as C^∞ , indicating that the functions are infinitely often differentiable. We define a *smooth structure* on \mathbb{M} to be a maximal collection of compatible coordinate charts, and call \mathbb{M} *smooth* if it has a smooth structure. A continuous function $f : \mathbb{M} \rightarrow \mathbb{R}$ is *smooth* if for each coordinate chart (U, ϕ) the composition $f \circ \phi^{-1}$ is smooth. A mapping $f : \mathbb{M} \rightarrow \mathbb{R}^3$ is *smooth* if all component functions $f_i = \pi_i \circ f$ are smooth, where π_i denotes projection onto the i -th factor.

For the time being, we assume that \mathbb{M} and f are smooth. If we choose a coordinate chart, we get a local parametrization of \mathbb{M} with two variables, s_1 and s_2 . Collecting the gradients of the coordinate functions in a matrix, we get the *Jacobian* of f :

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial s_1} & \frac{\partial f_1}{\partial s_2} \\ \frac{\partial f_2}{\partial s_1} & \frac{\partial f_2}{\partial s_2} \\ \frac{\partial f_3}{\partial s_1} & \frac{\partial f_3}{\partial s_2} \end{bmatrix}.$$

While this Jacobian matrix depends on the choice of local coordinates, its rank does not. Notice that the rank of the Jacobian is at most two. The mapping f is an *immersion* if the Jacobian has full rank two at all points of \mathbb{M} . It is an *embedding* if f is a homeomorphism onto its image. An embedding is necessarily an immersion, but not vice versa. For smooth mappings, there are three types of generic self-intersections, all illustrated in Figure II.8. The most interesting of the three is the

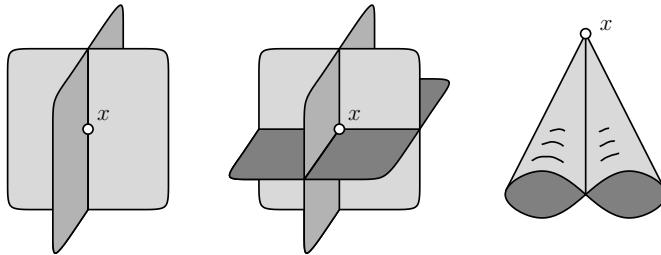


Figure II.8: From left to right: a double point, a triple point, a branch point.

branch point, which comes in several guises. We can construct it by cutting a disk from two sides toward the center, folding it, and re-gluing the sides as shown in Figure II.9. Embeddings have no self-intersections at all, and immersions have only the first two types and no branch points.

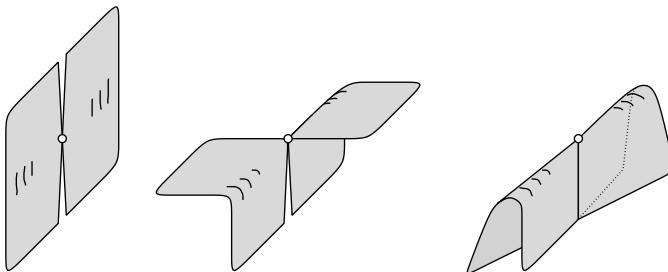


Figure II.9: Constructing the Whitney umbrella from a disk.

The piecewise linear case. The classification of generic self-intersections is similar in the piecewise linear case in which \mathbb{M} is given by a finite triangulation, K . However, in contrast to the smooth case, the enumeration of the generic types is elementary. Since \mathbb{M} is a 2-manifold, the triangles that contain a vertex form a disk. It is not difficult to see that imposing this condition on the vertices suffices to guarantee that K triangulates a 2-manifold without boundary. On the other hand, requiring that each edge belong to exactly two triangles is not sufficient.

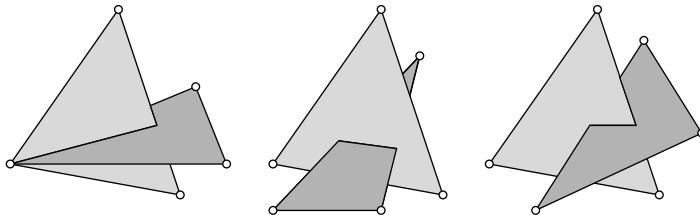


Figure II.10: The three ways that two triangles whose vertices are in general position in \mathbb{R}^3 can cross each other.

We put K into space by mapping each vertex to a point in \mathbb{R}^3 . The edges and triangles are mapped to the convex hulls of the images of their vertices. This mapping is an embedding iff any two triangles are either disjoint or they share a vertex or they share an edge. Any other type of intersection is improper and is referred to as a *crossing*. It is convenient to assume that the points are in general position, that is, no three are collinear and no four are coplanar. Under this assumption, there are only three types of crossings possible between two triangles, all shown in Figure II.10. Each crossing is a line segment common to two triangles. In the first case, one of the endpoints of the line segment coincides with the image of a vertex, which necessarily belongs to both crossing triangles. In the other two cases, each endpoint of the line segment lies on the images of an edge in the triangulation.

Recognizing crossings. We reduce the recognition problem from two triangles to an edge and a triangle and further to four points in space. Writing a_1, a_2, a_3 for the coordinates of the point a in space and similarly for the points x, y , and z , we say the sequence $axyz$ has *positive orientation* if the matrix

$$\Delta(a, x, y, z) = \begin{bmatrix} 1 & a_1 & a_2 & a_3 \\ 1 & x_1 & x_2 & x_3 \\ 1 & y_1 & y_2 & y_3 \\ 1 & z_1 & z_2 & z_3 \end{bmatrix}$$

has positive determinant. We observe that this corresponds to the case in which a sees xyz make a right turn in space. The four points lie in a common plane iff the determinant vanishes. Finally, we say $axyz$ has *negative orientation* if $\det \Delta(a, x, y, z) < 0$.

Using the ability to decide the orientation of a sequence of four points, we now turn to the next more complicated problem given by five points, a, b, x, y, z in \mathbb{R}^3 .

We say the edge ab *stabs* the triangle xyz if the two have an improper intersection. Assuming the five points are distinct and in general position, we have only two cases, namely either the intersection is empty or there is a point in the common interior of the edge and the triangle. Thus, ab stabs xyz iff a and b lie on different sides of the plane spanned by xyz and ab forms the same orientation with the three directed edges xy , yz , and zx .

```
boolean DOESSTAB(a, b, x, y, z)
    return sign(det Δ(a, x, y, z)) ≠ sign(det Δ(b, x, y, z)) and
        sign(det Δ(a, b, x, y)) = sign(det Δ(a, b, y, z)) = sign(det Δ(a, b, z, x)).
```

We finally return to the original recognition problem formulated for two triangles, abc and xyz . First, we consider the case in which they share one of the points, $a = x$. Then we have a crossing iff one of the respective opposite edges stabs the other triangle. Second, we consider the case in which the six points are distinct. Then the triangles are disjoint iff none of the six edges stabs the other triangle, and the triangles cross iff exactly two edges stab the other triangle. Assuming general position, there are no other cases. If the two stabbing edges belong to the same triangle, we have the case in the middle of Figure II.10, and if they belong to different triangles, we have the case on the right.

Curves and preimages. Returning to the case on the left in Figure II.10, we see that one endpoint of the line segment lies on the image of an edge in the triangulation; that is, it is not a vertex. There is a unique triangle on the other side of that edge that continues the intersections. Similarly, there are unique continuations of the intersection in the middle case and the right case. Starting at a crossing, we can therefore trace the intersection triangle by triangle, adding a line segment at a time. Since we only have finitely many triangles, the curve must either end or close up by coming back to where it started. These are the only two possibilities:

- a path that starts at the image of a vertex and ends at the image of another vertex;
- a closed curve that avoids the images of all vertices in the triangulation.

The first possibility involves two instances of the left case in Figure II.10; the second involves none. Almost all points of such a path or closed curve are double points. Exceptions are triple points at which the curves intersect each other or themselves. The number of triple points is at most the number of ways we can choose three triangles, which is finite, and generically there are no points that belong to more than three triangles.

When we trace a path or a closed curve in space, we can, at the same time, trace its preimage under the mapping f . In the case of a path, we get two arcs starting at a common vertex and ending at another common vertex of the triangulation. In the case of the closed curve, we get either two loops or one loop whose image covers the curve twice. The three cases are illustrated in Figure II.11. The most interesting

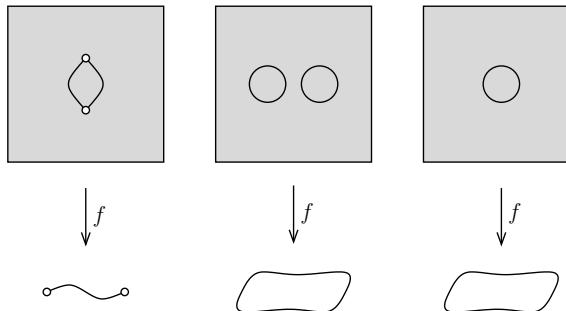


Figure II.11: The preimage of an intersection curve. From left to right: two arcs with common endpoints, two loops, one loop covering the closed curve twice.

case is the double-covering loop. Such a loop is necessarily orientation-reversing. To see this, we may again trace the closed curve, its image in \mathbb{R}^3 , and this time draw parallel curves to the left and the right on one of the two intersecting sheets. At the time we come back to where we started, the parallel curves have moved to the other sheet. There is either a clockwise or a counterclockwise rotation of the first sheet to the second that maps each curve locally to itself. If the rotation is clockwise, as seen by looking in the direction of the curve, then it is clockwise at all points of the curve. The same is true for the counterclockwise rotation. This implies that after another round we map the first sheet to itself but with reversed orientation. The double-covering loop can thus only happen if M is non-orientable. No conclusion can be drawn if the preimage consists of two loops.

To construct an example of a double-covering loop, we sweep the midpoint of a rod (a line segment) along a circle in space. The rod is normal to the circle at all times, but it may rotate within the normal plane as we sweep along. If there is no rotation, then the rod sweeps out a cylinder, and if the rotation is π after one time around, then we get a Möbius strip. However, if the rotation is $\frac{\pi}{2}$, we need a second time around to complete the surface. We thus get a Möbius strip that crosses itself along the center circle, which is covered twice.

Immersions of the Klein bottle. We have seen a first picture of the Klein bottle in Figure II.2. The surface in that drawing intersects itself along a path which ends at two branch points. In the smooth case, we get rank-deficient Jacobians at the branch points, implying that this is not the image of an immersion. However, the Klein bottle can also be mapped without branch points, and we conclude this section with the description of two such mappings.

In the first immersion, the neck of the bottle extends and turns back to the body, like a sleeping Flamingo, but then continues and passes through the surface, as sketched in Figure II.12 on the left. The closed intersection curve is the common image of two orientation-preserving loops. The second immersion is obtained by sweeping the cross point of a figure-eight curve along a circle in space. Similar to

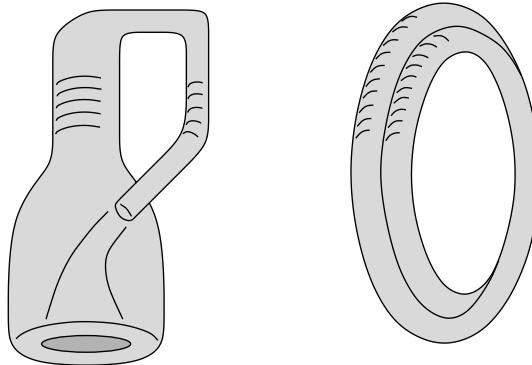


Figure II.12: Two immersions of the Klein bottle. Both models intersect themselves in a closed curve whose preimage is two loops. On the left, these loops are orientation-preserving, and on the right, they are orientation-reversing.

the rod example above, we keep the figure-eight normal to the circle at all times, but we rotate within the normal plane. Turning the figure-eight upside down during one time around, we exchange the lobes and form a surface that intersects itself along the circle, as sketched in Figure II.12 on the right. The preimage of the circle consists of two loops, both of which are orientation-reversing.

Bibliographic notes. The way surfaces in 3-dimensional space intersect each other and themselves is discussed in length and with many illustrations by Carter [28]. In the generic case, a smooth mapping to \mathbb{R}^3 has only three types of singularities: double points, triple points, and branch points [13]. Whitney proved that every d -manifold has an immersion in \mathbb{R}^{2d-1} [158]. This implies that every 2-manifold can be immersed in \mathbb{R}^3 . For the projective plane, we must have a branch point or a triple point, which implies that every immersion has a triple point [13]. Whitney also proved that every d -manifold can be embedded in \mathbb{R}^{2d} [157], implying that every 2-manifold can be embedded in \mathbb{R}^4 .

II.4 Surface Simplification

In applications, it is often necessary to simplify the data or its representation. One reason is measurement noise, which we would like to eliminate. Another reason is features, which we look for at various levels of resolution. In this section, we study edge contractions used in simplifying triangulated surface models of solid shapes.

Edge contraction. Suppose K is a triangulation of a 2-manifold without boundary. We recall that this means that edges are shared by pairs of triangles and vertices by rings of triangles, as depicted in Figure II.13. Let a and b be two vertices and ab the connecting edge in K . By the *contraction* of ab we mean the operation that

identifies a with b and removes duplicates from the triangulation. Calling the new vertex c , we get the new triangulation L from K by

- removing ab , abx , and aby ;
- substituting c for a and for b wherever they occur in the remaining set of vertices, edges, and triangles;
- removing resulting duplications making sure L is a set.

As a consequence of the operation, there are new incidences between edges and triangles that did not exist in K ; see Figure II.13.

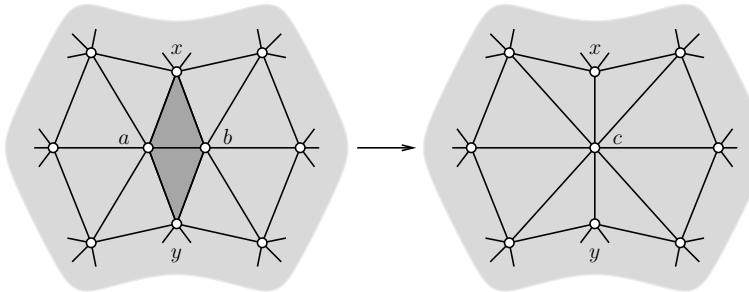


Figure II.13: To contract ab , we remove the two dark triangles and repair the hole by gluing their two left edges to their two right edges.

Algorithm. To simplify a triangulation, we iterate the edge contraction operation. In the abstract setting, any edge is as good as any other. In a practical situation, we will want to prioritize the edges so that contractions that preserve the shape of the manifold are preferred. To give meaning to this statement, we will define shape to mean the topological type of the surface as well as the geometric form we get when we embed the triangulation in \mathbb{R}^3 . We will discuss the latter meaning later and for now assume we have a function that assigns to each edge ab a non-negative real number $\text{ERROR}(ab)$ assessing the damage the contraction of ab causes to the geometric form. Small numbers will mean little damage. To write the algorithm, we assume a priority queue storing all edges ordered by the mentioned numerical error assessment. This is a data structure that supports the operations of returning the top priority edge as well as of inserting and deleting an edge, each in time at most logarithmic in the number of edges in the queue. Specifically, we assume a function `ISEMPTY` that tests whether or not the priority queue still contains edges and a function `MINEXTRACT` that removes the edge with minimum error from the priority queue and returns it. Furthermore, we assume the availability of a boolean test `ISSAFE` that decides whether or not the contraction of an edge preserves the topological type of the surface.

```

while not ISEMPTY do  $ab = \text{MINEXTRACT};$ 
    if ISSAFE( $ab$ ) then contract  $ab$  endif
endwhile.

```

Some modifications are necessary to recognize edges that no longer belong to the triangulation and to put edges back into the priority queue when they become safe for contraction. Details are omitted. The running time of the algorithm depends on the size of local neighborhoods in the triangulation and on the data structure we maintain to represent it. Under reasonable assumptions, the most time-consuming step is the maintenance of the priority queue, which for each step is only logarithmic in the number of edges.

Topological type. We consider the question of whether or not the contraction of an edge preserves the topological type. Define the *link* of an edge ab as the set of vertices that span triangles with ab , and define the link of a vertex a as the set of vertices that span edges with a and the set of edges that span triangles with a :

$$\begin{aligned} \text{Lk } ab &= \{x \in K \mid abx \in K\}; \\ \text{Lk } a &= \{x, xy \in K \mid ax, axy \in K\}. \end{aligned}$$

Since the topological type of K is that of a 2-manifold without boundary, each edge link is a pair of vertices and each vertex link is a closed curve made up of edges and vertices in K . Let L be obtained from K by contracting the edge ab . We show that the contraction of the edge ab preserves the topological type of the surface iff the links of the endpoints, a and b , meet in exactly two points, namely in the vertices x and y in the link of ab , as in Figure II.13. We will simplify the language by blurring the difference between a triangulation and the topological space it triangulates.

LINK CONDITION LEMMA. The triangulations K and L have the same topological type iff $\text{Lk } ab = \text{Lk } a \cap \text{Lk } b$.

PROOF. We have $\text{Lk } ab \subseteq \text{Lk } a, \text{Lk } b$, by definition. The only possible violation to the link condition is therefore an extra edge or vertex in the intersection of the two vertex links. If $\text{Lk } a$ and $\text{Lk } b$ share an edge, then the contraction of ab creates an edge that belongs to three triangles, contradicting the fact that L triangulates a 2-manifold. Similarly, if the two vertex links share no edge but a vertex $z \notin \text{Lk } ab$, then the contraction of ab creates an edge cz that belongs to four triangles, again contradicting the fact that L triangulates a 2-manifold.

To prove the other direction, we draw the link of c in L as a convex polygon in \mathbb{R}^2 ; see Figure II.14. Using Tutte's Theorem from the previous chapter, we can decompose the polygon by drawing the triangles incident to c in L . Similarly, we can decompose the polygon by drawing the triangles incident to a and b in K . We superimpose the two triangulations and refine to get a new triangulation, if necessary. The result is mapped back to K and to L , effectively refining the neighborhoods of a and b in K and that of c in L . The link of c and everything outside that link is untouched by the contraction. Hence, on and outside the link, K

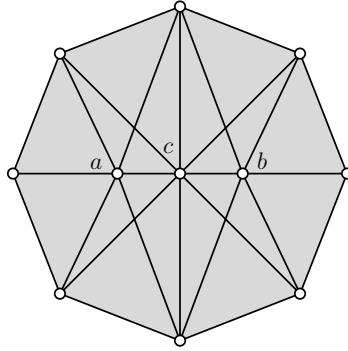


Figure II.14: Mapping the neighborhood of c in L to a triangulated polygon in the plane and overlaying it with a similar mapping of the neighborhoods of a and b in K .

and L are the same, and inside the link, K and L are now isomorphic by refinement. It follows that K and L are isomorphic and therefore have the same topological type. \square

Squared distance. To discuss the geometric meaning of shape, we now assume that K is embedded in \mathbb{R}^3 , with straight edges and flat triangles. To develop an error measure, we use the planes spanned by the triangles. Letting $u \in \mathbb{S}^2$ be the unit normal of a plane h and $\delta \in \mathbb{R}$ its offset, we can write h as the set of points $y \in \mathbb{R}^3$ for which $\langle y, u \rangle = -\delta$. Using matrix notation for the scalar product, the *squared distance* of a point $x \in \mathbb{R}^3$ from h is

$$d(x, h) = (x - y)^T \cdot u = x^T \cdot u + \delta,$$

where y is any point in the plane. Defining $\mathbf{x}^T = (x^T, 1)$ and $\mathbf{u}^T = (u^T, \delta)$, we can write this as a 4-dimensional scalar product, $\mathbf{x}^T \cdot \mathbf{u}$. We use this to express the sum of squared distances from a set of planes in matrix form. Letting H be a finite set of planes, this gives a function $E_H : \mathbb{R}^3 \rightarrow \mathbb{R}$ defined by

$$\begin{aligned} E_H(x) &= \sum_{h_i \in H} d^2(x, h_i) \\ &= \sum_{h_i \in H} (\mathbf{x}^T \cdot \mathbf{u}_i)(\mathbf{u}_i^T \cdot \mathbf{x}) \\ &= \mathbf{x}^T \cdot \left(\sum_{h_i \in H} \mathbf{u}_i \cdot \mathbf{u}_i^T \right) \cdot \mathbf{x}. \end{aligned}$$

Hence $E_H(x) = \mathbf{x}^T \cdot \mathbf{Q} \cdot \mathbf{x}$, where

$$\mathbf{Q} = \sum_{h_i \in H} (\mathbf{u}_i \cdot \mathbf{u}_i^T) = \begin{bmatrix} A & P & Q & U \\ P & B & R & V \\ Q & R & C & W \\ U & V & W & Z \end{bmatrix}$$

is a symmetric, four-by-four matrix that we refer to as the *fundamental quadric* of the map E_H . Writing $x^T = (x_1, x_2, x_3)$, we get

$$\begin{aligned} E_H(x) &= Ax_1^2 + Bx_2^2 + Cx_3^2 + 2(Px_1x_2 + Qx_1x_3 + Rx_2x_3) \\ &\quad + 2(Ux_1 + Vx_2 + Wx_3) + Z. \end{aligned}$$

We see that E_H is a quadratic map that is non-negative and unbounded.

Error assessment. In the application, we are interested in measuring the damage to the geometric form caused by contracting the edge ab to the new vertex c . We think of the operation as a map between vertices, $\varphi : \text{Vert } K \rightarrow \text{Vert } L$, defined by $\varphi(a) = \varphi(b) = c$ and $\varphi(x) = x$ for all $x \neq a, b$. Letting K_0 be the initial triangulation, we obtain L by a sequence of edge contractions giving rise to a composition of vertex maps, which is again a vertex map, $\varphi_0 : \text{Vert } K_0 \rightarrow \text{Vert } L$. The vertices in $V_c = \varphi_0^{-1}(c) \subseteq \text{Vert } K_0$ all map to c , and we let H be the set of planes spanned by triangles in K_0 incident to at least one vertex in V_c . Finally, we define the *error* of the contraction of ab as the minimum, over all possible placements of c as a point in \mathbb{R}^3 , of the sum of squared distances from the planes:

$$\text{ERROR}(ab) = \min_{c \in \mathbb{R}^3} E_H(c).$$

For generic sets of planes, this minimum is unique and easy to compute. The gradient of $E = E_H$ at a point x is the vector of steepest increase, $\nabla E(x) = (\frac{\partial E}{\partial x_1}(x), \frac{\partial E}{\partial x_2}(x), \frac{\partial E}{\partial x_3}(x))$. It is zero iff x minimizes E . The derivative with respect to x_i can be computed using the multiplication rule

$$\begin{aligned} \frac{\partial E}{\partial x_i} &= \frac{\partial \mathbf{x}^T}{\partial x_i} \cdot \mathbf{Q} \cdot \mathbf{x} + \mathbf{x}^T \cdot \mathbf{Q} \cdot \frac{\partial \mathbf{x}}{\partial x_i} \\ &= \mathbf{Q}[i]^T \cdot \mathbf{x} + \mathbf{x}^T \cdot \mathbf{Q}[i], \end{aligned}$$

where $\mathbf{Q}[i]$ is the i -th column and $\mathbf{Q}[i]^T$ is the i -th row of \mathbf{Q} . The point $c \in \mathbb{R}^3$ that minimizes E can thus be computed by setting $\frac{\partial E}{\partial x_i}$ to zero, for $i = 1, 2, 3$, and solving the resulting system of three linear equations.

Maintenance of the error measure. It can be expensive to compute the fundamental quadric from scratch but relatively inexpensive to maintain it throughout the algorithm. When we contract an edge ab , we associate the new vertex with the union of the two plane sets, $H_c = H_a \cup H_b$. Unfortunately, this is not a disjoint union, and we cannot just add the two quadrics. Instead, we use inclusion-exclusion and subtract the quadric of $H_{ab} = H_a \cap H_b$, which we store with the contracted edge. We describe how this works from the beginning.

Starting with the initial complex, K_0 , we store a quadric with every vertex, every edge, and every triangle. For a triangle abx , we store the quadric \mathbf{Q}_{abx} defined by the one plane that contains the triangle. An edge, ab , is shared by two triangles, abx and aby , and we store the quadric defined by the two corresponding planes, $\mathbf{Q}_{ab} = \mathbf{Q}_{abx} + \mathbf{Q}_{aby}$. A vertex, a , is shared by the ring of triangles in its star, and

we initialize its quadric, \mathbf{Q}_a , to the sum of the quadrics of these triangles. Note that the triangles that share the edge ab are precisely the ones that share both endpoints, a and b . This gives rise to a simple relationship between the sets of planes.

INVARIANT. Let abx be a triangle in the surface triangulation, with edges ab , ax , ay and vertices a , b , x . Then $H_{ab} = H_a \cap H_b$ and $H_{abx} = H_{ax} \cap H_{bx}$.

To maintain these two relations past an edge contraction, it is important that we limit ourselves to those that satisfy the Link Condition Lemma and therefore the topological type of the surface. The relations are therefore indeed invariants of the algorithm. Now consider the contraction of the edge ab . By the Invariant, the set of planes associated with the edge is the intersection of those of the endpoints. Hence we can compute the quadric of the new vertex as $\mathbf{Q}_c = \mathbf{Q}_a + \mathbf{Q}_b - \mathbf{Q}_{ab}$. We also get two new edges, cx and cy , and to maintain the Invariant, we associate each with the union of plane sets of the corresponding old edges. By the Invariant, these two sets overlap in the plane set of the shared triangle, which consists of a single plane. Hence, we get $\mathbf{Q}_{cx} = \mathbf{Q}_{ax} + \mathbf{Q}_{bx} - \mathbf{Q}_{abx}$ and $\mathbf{Q}_{cy} = \mathbf{Q}_{ay} + \mathbf{Q}_{by} - \mathbf{Q}_{aby}$.

Bibliographic notes. The algorithm described in this section is essentially the surface simplification algorithm by Garland and Heckbert [75]. They combine edge contractions with the error measure remembering the original form through accumulated quadrics. However, instead of maintaining the quadric through inclusion-exclusion, they take a short-cut and compute the quadric of the new vertex as the sum of quadrics of the endpoints of the contracted edge, without removing duplicates. In practice, this makes little difference because planes contribute at most in triplicates. The test for maintaining the topological type has been added later and more general versions of the Link Condition Lemma can be found in [48]. Priority queues are standard tools in computer science, and implementations are described in most texts on algorithms, including the third volume of Knuth's pioneering series on computer programming [94].

Exercises

The credit assignment reflects a subjective assessment of difficulty. A typical question can be answered using knowledge of the material combined with some thought and analysis.

1. **Classifying 2-manifolds** (two credits). Characterize the two surfaces depicted in Figure II.15 in terms of genus, boundary, and orientability.
2. **2-coloring** (two credits). Let K be a triangulation of an orientable 2-manifold without boundary. Construct L by decomposing each edge into two edges and each triangle into six triangles. To do this, we add a new vertex in the interior of each edge. Similarly, we add a new vertex in the interior of each triangle, connecting it to the six vertices in the boundary of the triangle. The resulting

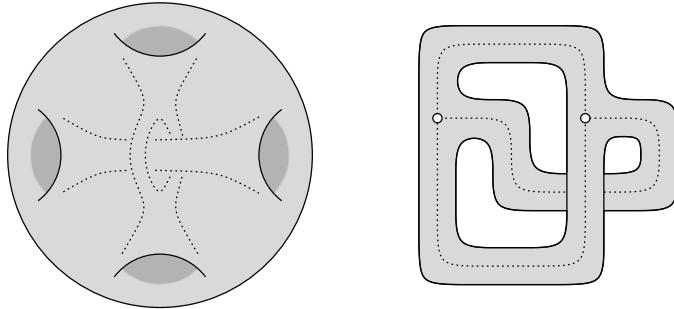


Figure II.15: Left: a 2-manifold without boundary obtained by adding tunnels inside the sphere. We see four tunnel openings and one tunnel passing though a fork of the other. Right: a 2-manifold with boundary obtained by thickening a graph.

structure is the same as the barycentric subdivision of K , which we will define in Chapter III.

- (i) Show that the vertices of L can be 3-colored such that no two neighboring vertices receive the same color.
- (ii) Prove that the triangles of L can be 2-colored such that no two triangles sharing an edge receive the same color.
3. **Klein bottle** (two credits). Cut and paste the standard polygonal schema for the Klein bottle (a, a, b, b) to obtain the polygonal schema in which opposite edges of a square are identified (a, b, a^{-1}, b) ; see Figure II.3.
4. **Triangulation of a 2-manifold** (two credits). Let $V = \{1, 2, \dots, n\}$ be a set of n vertices and $F \subseteq \binom{V}{3}$ a set of $\ell = \text{card } F$ triangles. Give an algorithm that takes time at most proportional to $n + \ell$ for the following tasks:
 - (i) decide whether or not every edge is shared by exactly two triangles;
 - (ii) decide whether or not every vertex belongs to a set of triangles whose union is a disk.
5. **Intersection tests in \mathbb{R}^3** (two credits). Let $a, b, c \in \mathbb{R}^3$ and $u, v, w \in \mathbb{R}^3$ be the vertices of two triangles in space. Write numerical tests for the following questions:
 - (i) Does u see a, b, c form a left turn or a right turn?
 - (ii) Does the line segment with endpoints u and v cross the plane that passes through a, b, c ?
 - (iii) Are the boundaries of the two triangles linked in \mathbb{R}^3 ?
6. **Irreducible triangulations** (two credits). An *irreducible* triangulation is one in which every edge contraction changes its topological type. Prove that the only irreducible triangulation of \mathbb{S}^2 is the boundary of the tetrahedron, which consists of four triangles sharing six edges and four vertices.

7. **Graphs on the Möbius strip** (one credit). Is every graph that can be embedded on the Möbius strip planar?
8. **Squared distance minimization** (two credits). Let S be a finite set of points in \mathbb{R}^3 and let $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ be defined by $f(x) = \sum_{p \in S} \|x - p\|^2$.
 - (i) Show that f is a quadratic function and has a unique minimum.
 - (ii) At which point does f attain its minimum?

Chapter III

Complexes

There are many ways to represent a topological space, one being a decomposition into simple pieces. This decomposition qualifies to be called a complex if the pieces are topologically simple and their common intersections are lower-dimensional pieces of the same kind. Within these requirements, we still have a great deal of freedom. Particularly attractive are the extreme choices: a few complicated or many simple pieces. The former choice lends itself to hand calculations of topological invariants but also to the design of aesthetically pleasing shapes, such as car bodies and the like. The latter choice is preferred in computation and automation. Since we focus on computational aspects of topology, we favor the latter extreme choice, of which the simplicial complex is the prime example.

III.1 Simplicial Complexes

In this book, we use simplicial complexes as the prime data structure to represent topological spaces. In this section, we introduce them in their geometric as well as abstract forms. The main technical result is the existence of simplicial maps that approximate continuous maps arbitrarily closely.

Simplices. Let u_0, u_1, \dots, u_k be points in \mathbb{R}^d . A point $x = \sum_{i=0}^k \lambda_i u_i$, with each $\lambda_i \in \mathbb{R}$, is an *affine combination* of the u_i if the λ_i sum to 1. The *affine hull* is the set of affine combinations. It is a *k-plane* if the $k+1$ points are *affinely independent*, by which we mean that any two affine combinations, $x = \sum \lambda_i u_i$ and $y = \sum \mu_i u_i$, are the same iff $\lambda_i = \mu_i$ for all i . The $k+1$ points are affinely independent iff the k vectors $u_i - u_0$, for $1 \leq i \leq k$, are linearly independent. In \mathbb{R}^d we can have at most d linearly independent vectors and therefore at most $d+1$ affinely independent points.

An affine combination, $x = \sum \lambda_i u_i$, is a *convex combination* if all λ_i are non-negative. The *convex hull* is the set of convex combinations. A *k-simplex* is the

convex hull of $k + 1$ affinely independent points, $\sigma = \text{conv} \{u_0, u_1, \dots, u_k\}$. We sometimes say the u_i span σ . Its dimension is $\dim \sigma = k$. We use special names for the first few dimensions: *vertex* for 0-simplex, *edge* for 1-simplex, *triangle* for 2-simplex, and *tetrahedron* for 3-simplex; see Figure III.1. Any subset of affinely

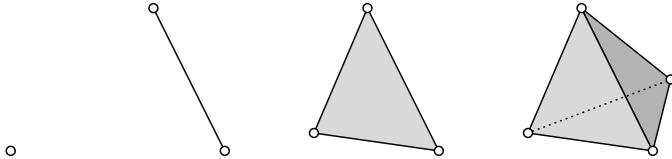


Figure III.1: From left to right: a vertex, an edge, a triangle, and a tetrahedron. We note that an edge has two vertices, a triangle has three edges, and a tetrahedron has four triangles as faces.

independent points is again affinely independent and therefore also defines a simplex. A *face* of σ is the convex hull of a non-empty subset of the u_i , and it is *proper* if the subset is not the entire set. We sometimes write $\tau \leq \sigma$ if τ is a face and $\tau < \sigma$ if it is a proper face of σ . If τ is a (proper) face of σ , we call σ a (*proper*) *coface* of τ . Since a set of size $k + 1$ has 2^{k+1} subsets, including the empty set, σ has $2^{k+1} - 1$ faces, all of which are proper except for σ itself. The *boundary* of σ , denoted as $\text{bd } \sigma$, is the union of all proper faces, and the *interior* is everything else, $\text{int } \sigma = \sigma - \text{bd } \sigma$. A point $x \in \sigma$ belongs to $\text{int } \sigma$ iff all its coefficients λ_i are positive. It follows that every point $x \in \sigma$ belongs to the interior of exactly one face, namely the one spanned by the points u_i that correspond to positive coefficients λ_i .

Simplicial complexes. We are interested in sets of simplices that are closed under taking faces and that have no improper intersections.

DEFINITION. A *simplicial complex* is a finite collection of simplices K such that $\sigma \in K$ and $\tau \leq \sigma$ implies $\tau \in K$, and $\sigma, \sigma_0 \in K$ implies $\sigma \cap \sigma_0$ is either empty or a face of both.

The *dimension* of K is the maximum dimension of any of its simplices. The *underlying space*, denoted as $|K|$, is the union of its simplices together with the topology inherited from the ambient Euclidean space in which the simplices live. A *polyhedron* is the underlying space of a simplicial complex. A *triangulation* of a topological space \mathbb{X} is a simplicial complex K together with a homeomorphism between \mathbb{X} and $|K|$. The topological space is *triangulable* if it has a triangulation. A *subcomplex* of K is a simplicial complex $L \subseteq K$. It is *full* if it contains all simplices in K spanned by vertices in L . A subcomplex of particular interest is the *j -skeleton* consisting of all simplices of dimension j or less, $K^{(j)} = \{\sigma \in K \mid \dim \sigma \leq j\}$. The 0-skeleton is also referred to as the *vertex set*, $\text{Vert } K = K^{(0)}$. Skeleta are generally not full.

A subset of a simplicial complex useful when discussing local neighborhoods is the *star* of a simplex τ consisting of all cofaces of τ , $\text{St } \tau = \{\sigma \in K \mid \tau \leq \sigma\}$. Generally, the star is not closed under taking faces. We can make it into a complex

by adding all missing faces. The result is the *closed star*, $\overline{\text{St}}\tau$, which is the smallest subcomplex that contains the star. The *link* consists of all simplices in the closed star that are disjoint from τ , $\text{Lk}\tau = \{v \in \overline{\text{St}}\tau \mid v \cap \tau = \emptyset\}$. If τ is a vertex, then the link is just the difference between the closed star and the star. More generally, it is the closed star minus the stars of all faces of τ . For example if K triangulates a 2-manifold without boundary, then the link of an edge is a pair of points, a 0-sphere, and the link of a vertex is a cycle of edges and vertices, a 1-sphere.

Abstract simplicial complex. It is often easier to construct a complex abstractly and worry about how to put it into Euclidean space later, if at all.

DEFINITION. An *abstract simplicial complex* is a finite collection of sets A such that $\alpha \in A$ and $\beta \subseteq \alpha$ implies $\beta \in A$.

The sets in A are its *simplices*. The *dimension* of a simplex is $\dim \alpha = \text{card } \alpha - 1$, and the dimension of the complex is the maximum dimension of any of its simplices. A *face* of α is a non-empty subset $\beta \subseteq \alpha$, which is *proper* if $\beta \neq \alpha$. The *vertex set* is the union of all simplices, $\text{Vert } A = \bigcup A$, that is, the set of all elements that lie in at least one simplex $\alpha \in A$. A *subcomplex* is an abstract simplicial complex $B \subseteq A$. Two abstract simplicial complexes are *isomorphic* if there is a bijection $b : \text{Vert } A \rightarrow \text{Vert } B$ such that $\alpha \in A$ iff $b(\alpha) \in B$. The largest abstract simplicial complex with a vertex set of size $n + 1$ is the n -dimensional simplex with a total number of $2^{n+1} - 1$ faces. Given a (geometric) simplicial complex K , we can construct an abstract simplicial complex A by throwing away all simplices and retaining only their sets of vertices. We call A a *vertex scheme* of K . Symmetrically, we call K a *geometric realization* of A . Constructing geometric realizations is surprisingly easy if the dimension of the ambient space is sufficiently high.

GEOMETRIC REALIZATION THEOREM. Every abstract simplicial complex of dimension d has a geometric realization in \mathbb{R}^{2d+1} .

PROOF. Let $f : \text{Vert } A \rightarrow \mathbb{R}^{2d+1}$ be an injection whose image is a set of points in general position. Specifically, any $2d + 2$ or fewer of the points are affinely independent. Let α and α_0 be simplices in A with $k = \dim \alpha$ and $k_0 = \dim \alpha_0$. The union of the two has size $\text{card}(\alpha \cup \alpha_0) = \text{card } \alpha + \text{card } \alpha_0 - \text{card}(\alpha \cap \alpha_0) \leq k + k_0 + 2 \leq 2d + 2$. The points in $\alpha \cup \alpha_0$ are therefore affinely independent, which implies that every convex combination x of points in $\alpha \cup \alpha_0$ is unique. Hence, x belongs to $\sigma = \text{conv } f(\alpha)$ as well as to $\sigma_0 = \text{conv } f(\alpha_0)$ iff x is a convex combination of $\alpha \cap \alpha_0$. This implies that the intersection of σ and σ_0 is either empty or the simplex $\text{conv } f(\alpha \cap \alpha_0)$, as required. \square

Simplicial maps. The natural counterparts of continuous maps between topological spaces are simplicial maps between simplicial complexes, which we now introduce. Let K be a simplicial complex with vertices u_0, u_1, \dots, u_n . Every point $x \in |K|$ belongs to the interior of exactly one simplex in K . For example, if $\sigma = \text{conv } \{u_0, u_1, \dots, u_k\}$ is this simplex, then we have $x = \sum_{i=0}^k \lambda_i u_i$ with

$\sum_{i=0}^k \lambda_i = 1$ and $\lambda_i > 0$ for all i . Setting $b_i(x) = \lambda_i$ for $0 \leq i \leq k$ and $b_i(x) = 0$ for $k+1 \leq i \leq n$, we have $x = \sum_{i=0}^n b_i(x)u_i$, and we call the $b_i(x)$ the *barycentric coordinates* of x in K .

We use these coordinates to construct a piecewise linear, continuous map, starting with a particular kind of map between the vertices of two simplicial complexes. A *vertex map* is a function $\varphi : \text{Vert } K \rightarrow \text{Vert } L$ with the property that the vertices of every simplex in K map to vertices of a simplex in L . Then φ can be extended to a continuous map $f : |K| \rightarrow |L|$ defined by

$$f(x) = \sum_{i=0}^n b_i(x)\varphi(u_i),$$

the *simplicial map* induced by φ . There is an alternative way to think of this construction. Fix a vertex u_j and consider the map $b_j : |K| \rightarrow \mathbb{R}$ which maps each point x to its j -th barycentric coordinate. The graph of this map has the shape of a hat, increasing from zero on and outside the link to one at u_j . The map b_j is continuous and is sometimes referred to as a basis function. The simplicial map is thus the weighted sum of the $n+1$ basis functions. To emphasize that the simplicial map is linear on every simplex, we usually drop the underlying space from the notation and write $f : K \rightarrow L$.

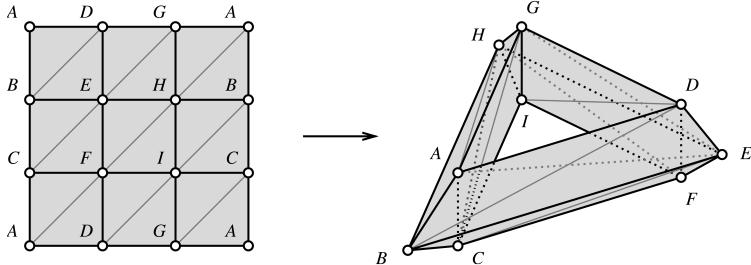


Figure III.2: A vertex map and its induced simplicial map from the square to the torus.

As an example, we consider the simplicial map $f : [0, 1]^2 \rightarrow \mathbb{T}^2$ illustrated in Figure III.2. Given the vertex map, the simplicial map is unique and glues the simplices of the triangulation of the square to obtain a triangulation of the torus. If the vertex map $\varphi : \text{Vert } K \rightarrow \text{Vert } L$ is bijective and $\varphi^{-1} : \text{Vert } L \rightarrow \text{Vert } K$ is also a vertex map, then the induced simplicial map f is a homeomorphism. In this case we call f a *simplicial homeomorphism* or an *isomorphism* between K and L .

Subdivisions. A simplicial complex L is a *subdivision* of another simplicial complex K if $|L| = |K|$ and every simplex in L is contained in a simplex in K . There are many ways to construct subdivisions. A particular one is the *barycentric subdivision*, $L = \text{Sd } K$, illustrated in Figure III.3. A crucial concept in its construction is the *barycenter* of a simplex, which is the average of its vertices. We proceed by

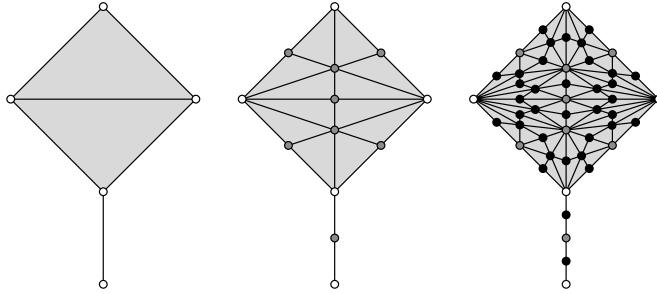


Figure III.3: Left: a simplicial complex consisting of two triangles, six edges, and five vertices. Middle and right: its first and second barycentric subdivisions.

induction over the dimension. To get started, the barycentric subdivision of the 0-skeleton is the same, $SdK^{(0)} = K^{(0)}$. Assuming we have the barycentric subdivision of $K^{(j-1)}$, we construct $SdK^{(j)}$ by adding the barycenter of every j -simplex as a new vertex and connecting it to the simplices that subdivide the boundary of the j -simplex.

The *diameter* of a set in Euclidean space is the supremum over the distances between its points. Since the simplices of K are point sets in Euclidean space, their diameters are well defined. The *mesh* of K is the maximum diameter of any simplex or, equivalently, the length of its longest edge.

MESH LEMMA. Letting δ be the mesh of the d -dimensional simplicial complex K , the mesh of SdK is at most $\frac{d}{d+1}\delta$.

PROOF. Let τ and v be complementary faces of a simplex $\sigma \in K$, that is, $\tau \cap v = \emptyset$ and $\dim \tau + \dim v = \dim \sigma - 1$. The line segment connecting the barycenters of τ and v has length at most δ . It splits into two edges in SdK at the barycenter of σ . Writing the barycenter of σ as a weighted sum of the barycenters of τ and v , we see that the lengths have proportions $1 + \dim v$ to $1 + \dim \tau$ which are both between $\frac{1}{k+1}$ and $\frac{k}{k+1}$, where $k = \dim \sigma$. It follows that both edges have length at most $\frac{k}{k+1} \leq \frac{d}{d+1}$ times δ . \square

By the Mesh Lemma, we can make the diameters of the simplices as small as we like by iterating the subdivision operation. For $n \geq 1$, define the n -th *barycentric subdivision* of K to be $Sd^n K = Sd(Sd^{n-1} K)$. As n goes to infinity, the mesh of $Sd^n K$ goes to zero.

Simplicial approximations. It is sometimes convenient to think of a vertex star as an open set of points. Formally, we define $N(u) = \bigcup_{\sigma \in St_u} \text{int } \sigma$. Let K and L be simplicial complexes. A continuous map $g : |K| \rightarrow |L|$ satisfies the *star condition* if the image of every vertex star in K is contained in a vertex star in L ; that is, for each vertex $u \in K$ there is a vertex $v \in L$ such that $g(N(u)) \subseteq N(v)$. Note that

we do not require, or even expect, v to be unique. Let $\varphi : \text{Vert } K \rightarrow \text{Vert } L$ map u to a vertex $\varphi(u) = v$ that exists by the star condition. To understand this new function, we take a point x in the interior of a simplex σ in K . Its image, $g(x)$, lies in the interior of a unique simplex τ in L . It follows that the star of every vertex u of σ maps into the star of a vertex v in L that contains the interior of τ . But this implies that v is a vertex of τ . We conclude that each vertex u of σ maps to a vertex $\varphi(u)$ of τ . Hence, φ is a vertex map and thus induces a simplicial map $f : K \rightarrow L$. This map satisfies the condition of a *simplicial approximation* of g , namely $g(N(u)) \subseteq N(f(u))$ for each vertex u of K .

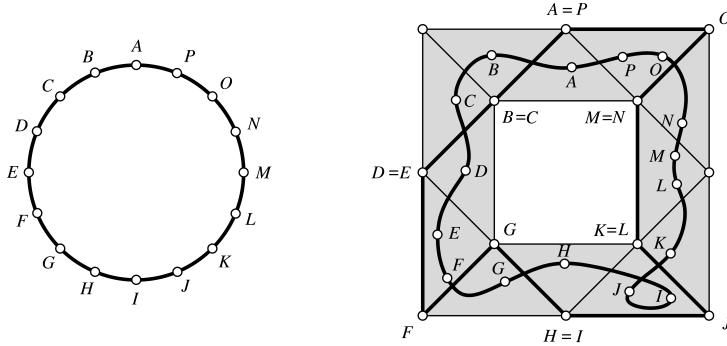


Figure III.4: The circle on the left is mapped into the closed annulus by a continuous map and a simplicial approximation of that map. Corresponding vertices are labeled by the same letter.

We illustrated the definitions in Figure III.4. The image we have in mind is that g and f are not too different. In particular, $g(x)$ and $f(x)$ belong to a common simplex in L for every $x \in |K|$. Given a continuous map $g : |K| \rightarrow |L|$, it is plausible that we can subdivide K sufficiently finely so that a simplicial approximation exists. To be sure, we prove this fact.

SIMPPLICIAL APPROXIMATION THEOREM. If $g : |K| \rightarrow |L|$ is continuous, then there is a sufficiently large integer n such that g has a simplicial approximation $f : \text{Sd}^n K \rightarrow L$.

PROOF. Cover $|K|$ with open sets of the form $g^{-1}(N(v))$, $v \in \text{Vert } L$. Since $|K|$ is compact, there is a positive real number λ such that any set of diameter less than λ is contained in one of the sets in the open cover. Choose n such that each simplex in $\text{Sd}^n K$ has diameter less than half of λ . Then each star in K has diameter less than λ , implying it lies in one of the sets $g^{-1}(N(v))$. Hence g satisfies the star condition, implying the existence of a simplicial approximation. \square

Bibliographic notes. The terminology we use for abstract and geometric simplicial complexes follows the one in Munkres [116]. The geometric realization of a d -dimensional abstract simplicial complex in \mathbb{R}^{2d+1} goes back to Karl Menger at the

beginning of the last century. We have seen that $2d + 1$ dimensions suffice for the geometric realization of any d -dimensional abstract simplicial complex. Complexes that require the full $2d + 1$ dimensions have been described by Flores [70] and van Kampen [143]. An example of such a complex is the d -skeleton of the $(2d + 2)$ -simplex, which does not embed in \mathbb{R}^{2d} . For $d = 1$ this is the complete graph of five vertices, which does not embed in the plane, as discussed in Chapter I.

A stronger version of the Simplicial Approximation Theorem played an important role in the development of combinatorial topology during the first half of the twentieth century. Known as the Hauptvermutung (German for “main conjecture”), it claimed that any two simplicial complexes that triangulate the same topological space have isomorphic subdivisions. This turned out to be correct for simplicial complexes of dimension 2 and 3 but not higher. The first counterexample found by Milnor was a simplicial complex of dimension 7 [110]. We refer to the book edited by Ranicki [125] for further information on the topic.

III.2 Convex Set Systems

Simplicial complexes often arise as intersection patterns of collections of sets. We begin with two fundamental results for convex sets and then proceed to the special case in which the sets are geometric balls.

Sets with common points. Let F be a finite collection of convex sets in \mathbb{R}^d . The smaller the dimension of the ambient Euclidean space, d , the more restrictive are the intersection patterns we observe. For example, if $d = 1$ and we have three intervals that intersect in pairs, then it is not possible that they do not intersect as a triplet. This result generalizes to higher dimensions.

HELLY’S THEOREM. Let F be a finite collection of closed, convex sets in \mathbb{R}^d . Every $d + 1$ of the sets have a non-empty common intersection iff they all have a non-empty common intersection.

PROOF. We prove the non-obvious direction by induction over the dimension, d , and the number of sets, $n = \text{card } F$. The implication is clearly true for $d = 1$ and all n , as well as for $n = d + 1$. Now suppose we have a minimal counterexample consisting of $n > d + 1$ closed, convex sets in \mathbb{R}^d , which we denote as X_1, X_2, \dots, X_n . By minimality of the counterexample, the set $Y_n = \bigcap_{i=1}^{n-1} X_i$ is non-empty and disjoint from X_n . Because Y_n and X_n are both closed and convex, we can find a $(d - 1)$ -dimensional plane h that separates and is disjoint from both sets, as in Figure III.5. Let F' be the collection of sets $Z_i = X_i \cap h$, for $1 \leq i \leq n - 1$, each a non-empty, closed, convex set in \mathbb{R}^{d-1} . By assumption, any d of the first $n - 1$ sets X_i have a common intersection with X_n . It follows that the common intersection of the d sets contains points on both sides of h , implying that any d of the sets Z_i have a non-empty common intersection. By minimality of the counterexample, this implies

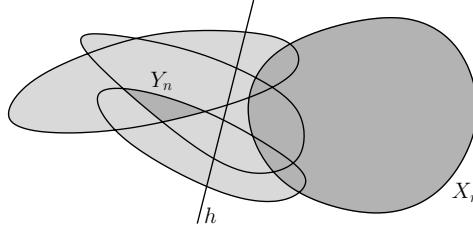


Figure III.5: The $(d - 1)$ -plane separates the n -th set from the common intersection of the first $n - 1$ sets in F .

$\bigcap F' \neq \emptyset$. This intersection is

$$\bigcap F' = \bigcap_{i=1}^{n-1} (X_i \cap h) = Y_n \cap h.$$

But this contradicts the choice of h as a $(d - 1)$ -plane disjoint from Y_n . \square

Convexity is a convenient but unnecessarily strong requirement in Helly's Theorem. Indeed, the conclusion holds if the sets in F are closed and all their non-empty common intersections are contractible, a property we will define shortly.

Homotopy type. We prepare the next step by introducing a notion of equivalence between topological spaces that is weaker than topological equivalence. We begin by considering two continuous maps, $f, g : \mathbb{X} \rightarrow \mathbb{Y}$. A *homotopy* between f and g is another continuous map $H : \mathbb{X} \times [0, 1] \rightarrow \mathbb{Y}$ that agrees with f for $t = 0$ and with g for $t = 1$; that is, $H(x, 0) = f(x)$ and $H(x, 1) = g(x)$ for all $x \in \mathbb{X}$. We may think of $t \in [0, 1]$ as time and the homotopy as a time-series of functions $f_t : \mathbb{X} \rightarrow \mathbb{Y}$ defined by $f_t(x) = H(x, t)$. It starts at $f_0 = f$ and ends at $f_1 = g$. Noting that this defines an equivalence relation, we write $f \simeq g$ and call f and g *homotopic* if there is a homotopy between them.

This notion can be used to relate spaces. Beginning with a special case, we call $\mathbb{Y} \subseteq \mathbb{X}$ a *retract* of \mathbb{X} if there is a continuous map $r : \mathbb{X} \rightarrow \mathbb{Y}$ with $r(y) = y$ for all $y \in \mathbb{Y}$. The map r is called a *retraction*. We call \mathbb{Y} a *deformation retract* and r a *deformation retraction* if there is a homotopy between r and the identity on \mathbb{X} , $r \simeq \text{id}_{\mathbb{X}}$. We also say that \mathbb{X} *deformation retracts* onto \mathbb{Y} . Clearly, every deformation retract is a retract but not the other way around. For example, a connected interval in the circle is a retract but not a deformation retract of \mathbb{S}^1 . We can also consider maps in both directions. Specifically, we call two not necessarily nested topological spaces, \mathbb{X} and \mathbb{Y} , *homotopy equivalent* if there are continuous maps $f : \mathbb{X} \rightarrow \mathbb{Y}$ and $g : \mathbb{Y} \rightarrow \mathbb{X}$ such that $g \circ f \simeq \text{id}_{\mathbb{X}}$ and $f \circ g \simeq \text{id}_{\mathbb{Y}}$. This gives an equivalence relation, and we write $\mathbb{X} \simeq \mathbb{Y}$ and say they have the same *homotopy type* if they are homotopy equivalent. The maps f and g are referred to as *homotopy equivalences* or *homotopy inverses* of each other.

To see that having the same homotopy type indeed generalizes being a deformation retract, we note that if $r : \mathbb{X} \rightarrow \mathbb{Y}$ is a deformation retraction, then $f = r$ and

g , the inclusion of \mathbb{Y} in \mathbb{X} , are continuous maps that satisfy the conditions and thus establish $\mathbb{X} \simeq \mathbb{Y}$. If \mathbb{Y} is a single point, then \mathbb{X} has the homotopy type of a point, and we say \mathbb{X} is *contractible*.

Nerves. We now return to our finite collection of sets, F . Without assuming the sets are convex, we define the *nerve* to consist of all non-empty subcollections whose sets have a non-empty common intersection:

$$\text{Nrv } F = \{X \subseteq F \mid \bigcap X \neq \emptyset\}.$$

It is always an abstract simplicial complex, no matter what sets we have in F . Indeed, if $\bigcap X \neq \emptyset$ and $Y \subseteq X$, then $\bigcap Y \neq \emptyset$. We can realize the nerve geometrically in some Euclidean space, so it makes sense to talk about its topology type and its homotopy type. We will sometimes do this without explicit construction of the geometric realization. As an example, consider the collection of four sets in Figure III.6 whose union is obviously not homotopy equivalent to the nerve. Nevertheless, taking the nerve preserves the homotopy type if the sets in the collection are convex. This is a fundamental result which we state formally but without proof.

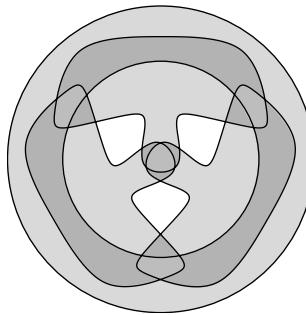


Figure III.6: A collection of four sets whose union is a disk with three holes in the plane. The nerve is the boundary complex of the tetrahedron which has the homotopy type of a sphere.

NERVE THEOREM. Let F be a finite collection of closed, convex sets in Euclidean space. Then the nerve of F and the union of the sets in F have the same homotopy type.

Similar to Helly's Theorem, the requirement on the sets can be relaxed without sacrificing the conclusion. Specifically, if $\bigcup F$ is triangulable, all sets in F are closed, and all non-empty common intersections are contractible, then $\text{Nrv } F \simeq \bigcup F$. We note that Helly's Theorem can be interpreted as a constraint on the structure of the nerve. Specifically, if the sets live in \mathbb{R}^d , then a subcollection of $k \geq d + 1$ sets cannot have all $\binom{k}{d+1}$ d -simplices in the nerve without having the entire k -simplex in the nerve.

Čech complexes. We now consider the special case in which the convex sets are closed geometric balls, all of the same radius, r . Let S be a finite set of points in \mathbb{R}^d and write $B_x(r) = x + r\mathbb{B}^d$ for the closed ball with center x and radius r . The Čech complex of S and r is the nerve of this collection of balls, but we substitute the center for each ball, that is,

$$\check{\text{C}}\text{ech}(r) = \{\sigma \subseteq S \mid \bigcap_{x \in \sigma} B_x(r) \neq \emptyset\}.$$

Clearly, a set of balls has a non-empty intersection iff their centers lie inside a common ball of the same radius. Indeed, a point y belongs to all balls iff $\|x - y\| \leq r$ for all centers x . An easy consequence of Helly's Theorem is therefore that every $d + 1$ points in S are contained in a common ball of radius r iff all points in S are. This is Jung's Theorem, which predates the more general theorem by Helly. The Čech complex does not necessarily have a geometric realization in \mathbb{R}^d , but it is fine as an abstract simplicial complex; see Figure III.7. For larger radius, the disks are

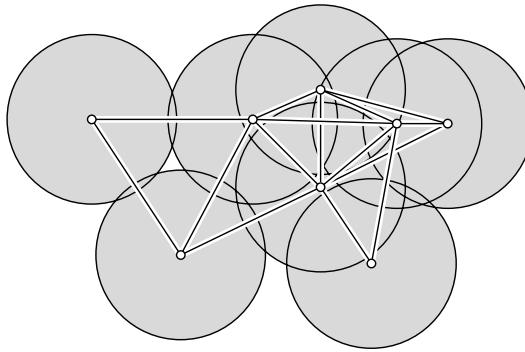


Figure III.7: Nine points with pairwise intersections among the disks indicated by straight edges connecting their centers. The Čech complex fills nine of the ten possible triangles as well as the two tetrahedra. The only difference between the Vietoris-Rips and the Čech complexes is the tenth triangle, which belongs only to the former.

bigger and create more overlaps while retaining the ones for smaller radius. Hence $\check{\text{C}}\text{ech}(r_0) \subseteq \check{\text{C}}\text{ech}(r)$ whenever $r_0 \leq r$. If we continuously increase the radius, from 0 to ∞ , we get a discrete family of nested Čech complexes. We will come back to this construction later.

Smallest enclosing balls. Beyond sets of two points, it seems cumbersome to recognize the ones that form simplices in the Čech complex. Nevertheless, there is a fast algorithm for this purpose.

Let $\sigma \subseteq S$ be a subset of the given points. We have seen that deciding whether or not σ belongs to $\check{\text{C}}\text{ech}(r)$ is equivalent to deciding whether or not σ fits inside a ball of radius r . Let the *miniball* of σ be the smallest closed ball that contains σ , which we note is unique. The radius of the miniball is smaller than or equal to

r iff $\sigma \in \check{\text{Cech}}(r)$, so finding it solves our problem. Observe that the miniball is already determined by a subset of $k + 1 \leq d + 1$ of the points, which all lie on its boundary. If we know this subset, then we can verify the miniball by testing that it indeed contains all the other points. In a situation in which we have many more points than dimensions, the chance that a point belongs to this subset is small and discarding it is easy. This is the strategy of the Miniball Algorithm. It takes two disjoint subsets τ and v of σ and returns the miniball that contains all points of τ in its interior and all points of v on its boundary. To get the miniball of σ , we call $\text{MINIBALL}(\sigma, \emptyset)$.

```

ball MINIBALL( $\tau, v$ )
  if  $\tau = \emptyset$  then compute the miniball  $B$  of  $v$  directly
    else choose a random point  $u \in \tau$ ;
       $B = \text{MINIBALL}(\tau - \{u\}, v)$ ;
      if  $u \notin B$  then
         $B = \text{MINIBALL}(\tau - \{u\}, v \cup \{u\})$ 
      endif
    endif; return  $B$ .
  
```

When τ is empty, we have a set v of at most $d + 1$ points, which we know all lie on the boundary. Assuming the dimension, d , is a constant, we can compute their miniball directly and in constant time. To analyze the running time, we ask how often we execute the test “ $u \notin B$ ”. Let $t_j(n)$ be the expected number of such tests for calling MINIBALL with n points in τ and $j = d + 1 - \text{card } v$ possibly open positions on the boundary of the miniball. Obviously, $t_j(0) = 0$, and it is reassuring that the constant amount of work needed to compute the ball for the at most $d + 1$ points in v is paid for by the test that initiated the call. Consider $n > 0$. We have one call with parameters $n - 1$ and j , one test “ $u \notin B$ ”, and one call with parameters $n - 1$ and $j - 1$. The probability that the second call indeed happens is at most j out of n . Hence,

$$t_j(n) \leq t_j(n - 1) + 1 + \frac{j}{n} \cdot t_{j-1}(n - 1).$$

Setting $j = 0$, we get $t_0(n) \leq t_0(n - 1) + 1$ and therefore $t_0(n) \leq n$. Similarly, $t_1(n) \leq t_1(n - 1) + 2 \leq 2n$. More generally, we get $t_j(n) \leq (j + 1)!n$, which is a constant times n since $j \leq d + 1$ is a constant. In summary, for constant dimension the algorithm takes constant time per point in the expected case.

Vietoris-Rips complexes. Instead of checking all subcollections, we may just check pairs and add 2- and higher-dimensional simplices whenever all their edges are in the complex. This simplification leads to the *Vietoris-Rips complex* of S and r consisting of all subsets of diameter at most $2r$:

$$\text{Vietoris-Rips}(r) = \{\sigma \subseteq S \mid \text{diam } \sigma \leq 2r\}.$$

Clearly, the edges in the Vietoris-Rips complex are the same as in the $\check{\text{Cech}}$ complex. Furthermore, $\check{\text{Cech}}(r) \subseteq \text{Vietoris-Rips}(r)$ because the latter contains every simplex

warranted by the given edges. We now prove that the containment relation can be reversed if we are willing to increase the radius in the definition of the Čech complex by a multiplicative constant.

VIETORIS-RIPS LEMMA. Letting S be a finite set of points in some Euclidean space and letting $r \geq 0$, we have $\text{Vietoris-Rips}(r) \subseteq \check{\text{Cech}}(\sqrt{2}r)$.

PROOF. A simplex is *regular* if all its edges have the same length. A convenient representation for dimension d is the *standard d -simplex*, Δ^d , spanned by the endpoints of the unit coordinate vectors in \mathbb{R}^{d+1} ; see Figure III.8. Each edge of Δ^d has

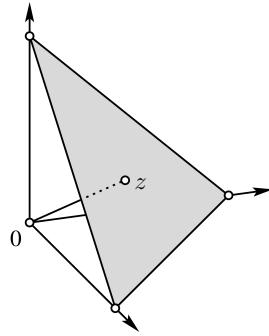


Figure III.8: The standard triangle connecting the unit coordinate vectors in \mathbb{R}^3 .

length $\sqrt{2}$. By symmetry, the distance of the origin from the standard simplex is its distance from the barycenter, the point z whose $d+1$ coordinates are all equal to $\frac{1}{d+1}$. That distance is therefore $\|z\| = 1/\sqrt{d+1}$. The barycenter is also the center of the smallest d -sphere that passes through the vertices of Δ^d . Writing r_d for the radius of that sphere, we have $r_d^2 = 1 - \|z\|^2 = \frac{d}{d+1}$. For dimension 1, this is indeed half the length of the interval, and for dimension 2, it is the radius of the equilateral triangle. As the dimension goes to infinity, the radius grows and approaches 1 from below. Any set of $d+1$ or fewer points for which the same d -ball of radius r_d is the miniball has a pair at distance $\sqrt{2}$ or larger. It follows that every simplex of diameter $\sqrt{2}$ or less belongs to $\check{\text{Cech}}(r_d)$. Multiplying with $\sqrt{2}r$, we get $\text{Vietoris-Rips}(r) \subseteq \check{\text{Cech}}(\sqrt{2}rr_d)$. Since $r_d \leq 1$ for all d , the latter is a subcomplex of $\check{\text{Cech}}(\sqrt{2}r)$, which implies the claimed subcomplex relationship. \square

Bibliographic notes. Helly proved his theorem at the beginning of the last century, first for convex sets and then for sets with contractible common intersections [83, 84]. The concept of nerve was introduced at about the same time by Alexandrov [9]. The Nerve Theorem goes back to Borsuk [19], Leray [102], and others. It has a complicated literature, with versions differing in the generality of the assumption and the strength of the conclusion. The Čech complexes are inspired by the theory of Čech homology, from which they borrow their name. The Vietoris-Rips complex appears in Vietoris [146] and in later work by Rips; see [79]. Algorithms for finding

the smallest ball enclosing a finite set of points have been studied in computational geometry, culminating in the randomized minidisk algorithm of Welzl which has versions that are efficient even for large sets in high dimensions [155].

III.3 Delaunay Complexes

In this section, we introduce a geometric constructions that limits the dimension of the simplices we get from a nerve. The main new structures are the Voronoi diagram and the Delaunay complex of a finite set of points. We begin by studying the inversion of space through a unit sphere.

Inversion. Recall that \mathbb{S}^d is the d -dimensional sphere with center at the origin and unit radius in \mathbb{R}^{d+1} . To invert \mathbb{R}^{d+1} , we map each point $x \neq 0$ to the point on the same half-line whose distance from the origin is the reciprocal of the distance of x from 0. More formally, the *inversion* maps x to $\iota(x) = x/\|x\|^2$. It exchanges inside with outside and leaves points on \mathbb{S}^d fixed. Clearly, $\iota(\iota(x)) = x$. We construct the image of a point x inside \mathbb{S}^d by drawing right-angled triangles. First, we find a point $p \in \mathbb{S}^d$ such that $0xp$ has a right angle at x . Second, we choose x' on the half-line of x such that $0px'$ has a right angle at p . The angle at 0 is the same in both so the two triangles are similar. Hence, $\|x\| : \|p\| = \|p\| : \|x'\|$, which implies $\|x\|\|x'\| = \|p\|^2 = 1$ and thus $x' = \iota(x)$. We use this construction to show that the inversion maps spheres to spheres. We note, however, that it generally does not map centers to centers.

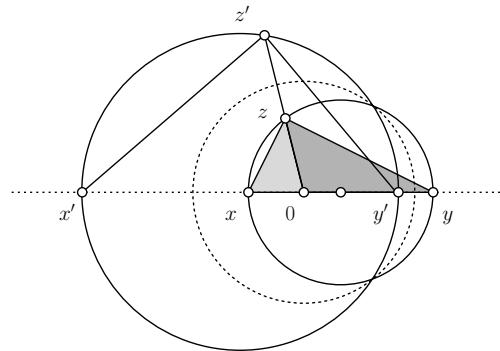


Figure III.9: The dotted circle represents the sphere \mathbb{S}^d centered at the origin. As z sweeps out the circle passing through x and y , its image, $z' = \iota(z)$, sweeps out the circle passing through x' and y' .

INVERSION LEMMA. Let Σ be a d -sphere in \mathbb{R}^{d+1} . If $0 \notin \Sigma$, then $\iota(\Sigma)$ is a d -sphere, and if $0 \in \Sigma$, then $\iota(\Sigma)$ is a d -plane.

PROOF. Consider first the case in which Σ does not pass through the origin, as in Figure III.9. If 0 is the center of Σ , then the result is obvious, so assume 0 is not the center. Draw the line passing through 0 and the center; it intersects Σ in points x and y , which we invert to get points $x' = \iota(x)$ and $y' = \iota(y)$. Let z be another point on Σ and let $z' = \iota(z)$ be its inverse. Then $\|x\|\|x'\| = \|z\|\|z'\| = 1$, which implies that the triangles $0xz$ and $0z'x'$ are similar. By the same token, $0yz$ and $0z'y'$ are similar. But xyz has a right angle at z , implying the angles at x' and y' inside $x'y'z'$ add up to a right angle. It follows that $x'y'z'$ has a right angle at z' . As z travels on Σ , the sphere with diameter xy , the image z' travels on $\iota(\Sigma)$, the sphere with diameter $x'y'$. What happens when Σ passes through the origin, say $0 = x$? Then the triangle $0y'z'$ has a right angle at y' . Equivalently, the image of Σ is the plane normal to the vector y and passing through the point y' . \square

The Inversion Lemma suggests that we think of a d -plane as a special kind of d -sphere, namely one that passes through the point at infinity.

Stereographic projection. Inversion can be defined relative to any center $z \in \mathbb{R}^{d+1}$ and any radius $r > 0$, that is, $\iota_{z,r}(x) = r \cdot \iota\left(\frac{x-z}{r}\right) + z$. It is not difficult to check that x and $x' = \iota_{z,r}(x)$ indeed lie on the same half-line emanating from z and the product of their distances is $\|x - z\|\|x' - z\| = r^2$, as desired. We consider

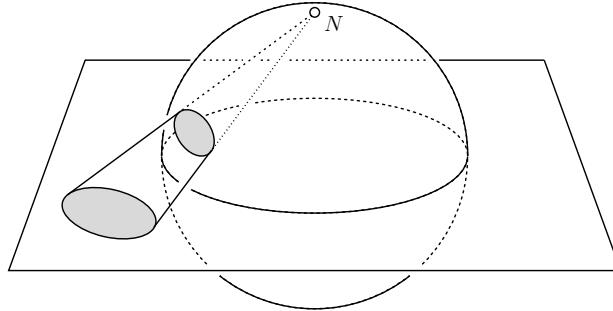


Figure III.10: The stereographic projection maps a circle on the unit sphere to a circle in the plane. If the circle on the sphere passes through the north pole then its image is a line, that is, a circle that passes through the point at infinity.

the special case in which the center is the point $N = (0, \dots, 0, 1)$, the north pole of \mathbb{S}^d , and the radius is $r = \sqrt{2}$, the Euclidean distance between the north pole and the equator. The image of \mathbb{S}^d is the d -plane of points with vanishing $(d+1)$ -st coordinates, which we denote as \mathbb{R}^d . The *stereographic projection* is the restriction of this particular inversion to the unit sphere, that is, $\varsigma : \mathbb{S}^d - \{N\} \rightarrow \mathbb{R}^d$ defined by $\varsigma(x) = \iota_{N,\sqrt{2}}(x)$, as sketched in Figure III.10. Similar to the inversion, the stereographic projection preserves spheres.

STEREOGRAPHIC PROJECTION LEMMA. Let Σ' be a $(d-1)$ -sphere on \mathbb{S}^d . If $N \notin \Sigma'$, then $\varsigma(\Sigma')$ is a $(d-1)$ -sphere, and if $N \in \Sigma'$, then $\varsigma(\Sigma')$ is a $(d-1)$ -plane.

Indeed, every $(d - 1)$ -sphere considered in the lemma is the intersection of \mathbb{S}^d with another d -sphere. Its image is therefore the intersection of \mathbb{R}^d with the image of the d -sphere, which is either a d -sphere or a d -plane. The intersection is thus either a $(d - 1)$ -sphere or a $(d - 1)$ -plane. As before, we consider a plane as a special sphere that passes through the point at infinity.

Voronoi diagram. We will use stereographic projection and the more general inversion to elucidate the construction of a particular simplicial complex, called the Delaunay complex, from a finite set $S \subseteq \mathbb{R}^d$. As a first step, we define the *Voronoi cell* of a point u in S as the set of points for which u is the closest, $V_u = \{x \in \mathbb{R}^d \mid \|x - u\| \leq \|x - v\|, v \in S\}$. It is the intersection of half-spaces of points at least

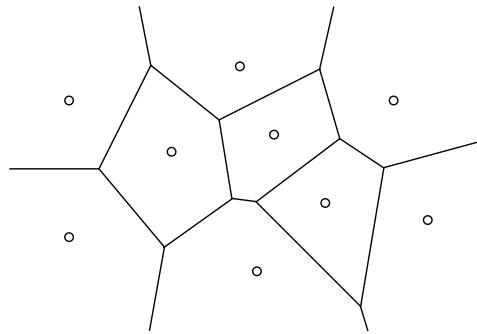


Figure III.11: The Voronoi diagram of nine points in the plane. By definition, each vertex of the diagram is equally far from the points that generate the incident Voronoi cells and further from all other points in S .

as close to u as to v , over all points v in S . Hence, V_u is a convex polyhedron in \mathbb{R}^d . Any two Voronoi cells meet at most in a common piece of their boundary, and together the Voronoi cells cover the entire space, as illustrated in Figure III.11. The *Voronoi diagram* of S is the collection of Voronoi cells of its points.

We will shortly use a generalization of the concept to points u with real weights w_u . The *weighted squared distance*, or *power*, of a point $x \in \mathbb{R}^d$ from u is $\pi_u(x) = \|x - u\|^2 - w_u$. For positive weight, we can interpret the weighted point as the sphere with center u and square radius w_u . For a point x outside this sphere, the power is positive and equal to the square length of a tangent line segment from x to the sphere. For x on the sphere the power vanishes, and for x inside the sphere the power is negative. The *bisector* of two weighted points is the set of points with equal power from both. Just as in the unweighted case, the bisector is a plane normal to the line connecting the two points, except that it is not necessarily halfway between them; see Figure III.12. Given a finite set of weighted points, we can thus define the *weighted Voronoi cell*, or *power cell*, of u as the set of points $x \in \mathbb{R}^d$ with $\pi_u(x) \leq \pi_v(x)$ for all weighted points v in the set. Finally, the *weighted Voronoi diagram*, or *power diagram*, is the set of power cells of the weighted points.

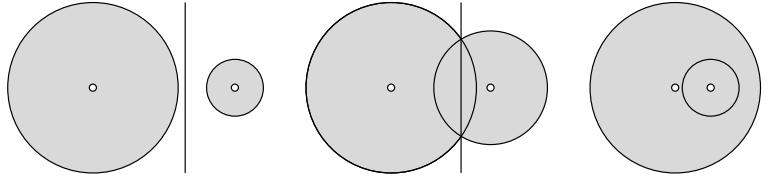


Figure III.12: The bisectors of pairs of weighted points. From left to right: two disjoint circles side by side, two intersecting circles, and two nested circles.

Lifting. We get a different and perhaps more illuminating view of the Voronoi diagram by lifting its cells to one higher dimension. Let S be a finite set of points in \mathbb{R}^d , as before, but draw them in \mathbb{R}^{d+1} , adding zeros as $(d+1)$ -st coordinates. Map each point $u \in S$ to \mathbb{S}^d using the inverse of stereographic projection, and let Π_u be the d -plane tangent to \mathbb{S}^d touching the sphere in the point $\varsigma^{-1}(u)$, as illustrated in Figure III.13. Using inversion, we now map each d -plane Π_u to the d -sphere

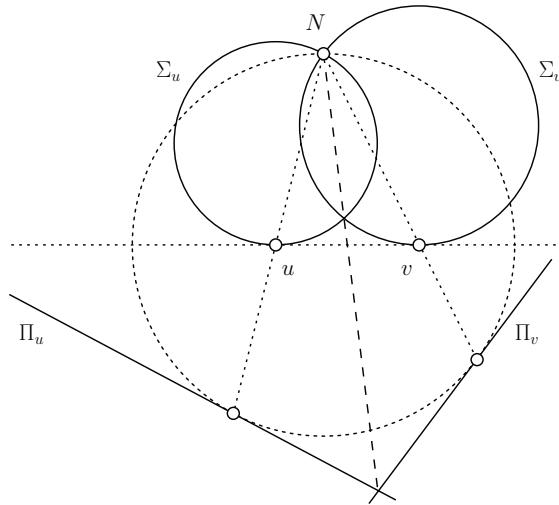


Figure III.13: We map the points u and v in \mathbb{R}^1 to the lines Π_u and Π_v tangent to \mathbb{S}^1 and further to the circles Σ_u and Σ_v passing through N and tangent to \mathbb{R}^1 . The dashed line connecting N and the midpoint between u and v passes through the intersection of the two circles and the intersection of the two lines.

$\Sigma_u = \iota(\Pi_u)$. It passes through the north pole and is tangent to \mathbb{R}^d , the preimage of \mathbb{S}^d . The arrangements of planes and of spheres are closely related to the Voronoi diagram. We focus on the spheres first.

FIRST SPHERE LEMMA. A point $x \in \mathbb{R}^d$ belongs to the Voronoi cell of $u \in S$ iff the first intersection of the directed line segment from x to N with the d -spheres defined by the points in S is with Σ_u .

PROOF. Interpret the sphere Σ_u as a weighted point, namely its center with weight equal to the square of its radius. The power of a point x is the squared length of a tangent line segment, which is equal to $\|x - u\|^2$ if $x \in \mathbb{R}^d$. It follows that the weighted Voronoi cell of the weighted center intersects \mathbb{R}^d in the Voronoi cell of u . The claim follows because all bisectors of the weighted points pass through N . \square

Switching from spheres to planes, we get a similar characterization of the Voronoi diagram in terms of tangent planes.

FIRST PLANE LEMMA. A point $x \in \mathbb{R}^d$ belongs to the Voronoi cell of $u \in S$ iff the first intersection of the directed line segment from N to x with the d -planes defined by the points in S is with Π_u .

Delaunay triangulation. The *Delaunay complex* of a finite set $S \subseteq \mathbb{R}^d$ is isomorphic to the nerve of the Voronoi diagram:

$$\text{Delaunay} = \{\sigma \subseteq S \mid \bigcap_{u \in \sigma} V_u \neq \emptyset\}.$$

We say the set S is in general position if no $d + 2$ of the points lie on a common $(d-1)$ -sphere. This assumption implies that no $d+2$ Voronoi cells have a non-empty common intersection. Equivalently, the dimension of any simplex in the Delaunay complex is at most d . Assuming general position, we get a geometric realization by taking convex hulls of abstract simplices, as in Figure III.14. The result is often referred to as the *Delaunay triangulation* of S . To see that this construction indeed

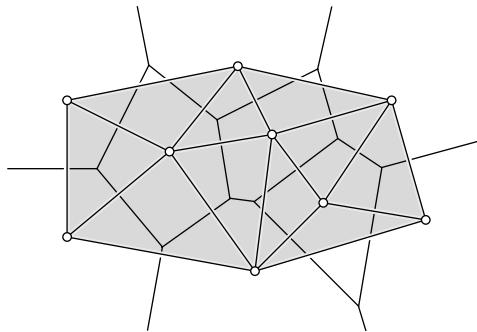


Figure III.14: The Delaunay triangulation superimposed on the Voronoi diagram. No four of the given points are cocircular implying the Delaunay complex has simplices of dimension at most two and a canonical geometric realization in \mathbb{R}^2 .

gives a geometric realization of the Delaunay complex, we lift the points to the set $\varsigma^{-1}(S)$ on \mathbb{S}^d . Similarly, we lift a general point $x \in \mathbb{R}^d$ to the d -plane Π_x tangent to \mathbb{S}^d at the point $\varsigma^{-1}(x)$. Keeping the same normal direction, we move this plane toward N . This corresponds to growing a $(d-1)$ -sphere around x . The first point encountered by the plane corresponds to the first point encountered by the sphere,

which is therefore the nearest to x . This suggests we add N to the set of lifted points and we take the convex hull in \mathbb{R}^{d+1} . The boundary of the resulting convex polytope consists of faces up to dimension d , some of which share N as a vertex. We are interested in the other faces, since they are spanned by points that correspond to Voronoi cells with a non-empty common intersection. Using the central projection from N , we map these faces to \mathbb{R}^d . By convexity of the polytope, the images of the faces have no improper intersections. Indeed, we get the geometric realization of the Delaunay complex, as promised.

Similar to the Voronoi diagram, we can generalize the Delaunay complex to a finite set of points with real weights. Specifically, the *weighted Delaunay complex* is the abstract simplicial complex that contains a subset of the weighted points iff their weighted Voronoi cells have a non-empty common intersection. In contrast to the unweighted case, the cell of a weighted point can be empty, a difference that is sometimes overlooked. As a consequence, the vertex set of the weighted Delaunay triangulation is a subset and not necessarily the entire set of given weighted points. Assuming general position, this complex can again be geometrically realized by taking convex hulls of the abstract simplices. The appropriate notion of general position is that no point of \mathbb{R}^d has the same power from more than $d + 1$ of the weighted points. This property is satisfied with probability one, a necessary requirement for a general position assumption.

Bibliographic notes. Voronoi diagrams are named after Georgy Voronoi [149, 150], and Delaunay triangulations are named after Boris Delaunay (also Delone) [44]. Both structures had been studied centuries earlier by others, including Dirichlet, Gauß, and Descartes. Weighted Voronoi diagrams are perhaps as old as the unweighted ones and are known under a plethora of different names, including Thiessen polygons, Dirichlet tessellations, and power diagrams; see the survey article by Aurenhammer [11]. Their dual weighted Delaunay triangulations are also known under a variety of names, including regular triangulations and coherent triangulations; see e.g. [76]. Algorithms for constructing Delaunay triangulations with an emphasis on mesh generation applications can be found in [54].

III.4 Alpha Complexes

In this section, we use a radius constraint to introduce a family of subcomplexes of the Delaunay complex. These complexes are similar to the Čech complexes but differ from them by having canonical geometric realizations.

Union of balls. Let S be a finite set of points in \mathbb{R}^d and r a non-negative real number. As before, for each $u \in S$, we let $B_u(r) = u + r\mathbb{B}^d$ be the closed ball with center u and radius r . The union of these balls is the set of points at distance at most r from at least one of the points in S . To decompose the union, we intersect each ball with the corresponding Voronoi cell, $R_u(r) = B_u(r) \cap V_u$. Since balls and Voronoi cells are convex, the $R_u(r)$ are also convex. Any two of them are disjoint

or overlap along a common piece of their boundaries, and together the $R_u(r)$ cover the entire union, as in Figure III.15. The *alpha complex* is the nerve of this cover, but we substitute the center for each ball, that is,

$$\text{Alpha}(r) = \{\sigma \subseteq S \mid \bigcap_{u \in \sigma} R_u(r) \neq \emptyset\}.$$

Since $R_u(r) \subseteq V_u$, the alpha complex is a subcomplex of the Delaunay complex. It follows that for a set S in general position, we get a geometric realization by taking convex hulls of abstract simplices, just as in the previous section. Furthermore, $R_u(r) \subseteq B_u(r)$, which implies $\text{Alpha}(r) \subseteq \check{\text{C}}\text{ech}(r)$. Since the $R_u(r)$ are closed and convex and together they cover the union, the Nerve Theorem implies that the union of balls and $\text{Alpha}(r)$ have the same homotopy type: $\bigcup_{u \in S} B_u(r) \simeq |\text{Alpha}(r)|$.

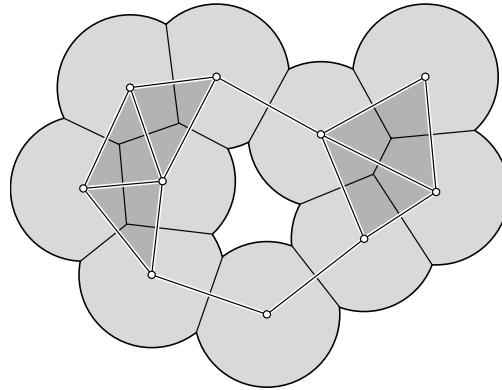


Figure III.15: The union of disks is decomposed into convex regions by the Voronoi diagram. The corresponding alpha complex is superimposed.

Weighted alpha complexes. For many applications, it is useful to permit balls with different sizes. An example of significant importance is the modeling of biomolecules, such as proteins, RNA, and DNA. Each atom is represented by a ball whose radius reflects the range of its van der Waals interactions and thus depends on the atom type. Therefore, let S be a finite set of points u with real weights w_u . As in the previous section, we think of u as a ball B_u with center u and squared radius $r_u^2 = w_u$. We again consider the union of the balls, which we decompose into convex regions, now using weighted Voronoi cells, $R_u = B_u \cap V_u$. This is illustrated in Figure III.16. In analogy to the unweighted case, the *weighted alpha complex* of S is the nerve of the collection of regions R_u , but we again substitute the points for the regions. Equivalently, it is the set of abstract simplices $\sigma \subseteq S$ such that $\bigcap_{u \in \sigma} R_u \neq \emptyset$. The weighted alpha complex is a subcomplex of the weighted Delaunay complex. Assuming the weighted points are in general position, we again get a geometric realization by taking convex hulls of abstract simplices. It will be convenient to blur the difference, which we do by using the exact same notation and dropping the term *weighted* unless it is essential.

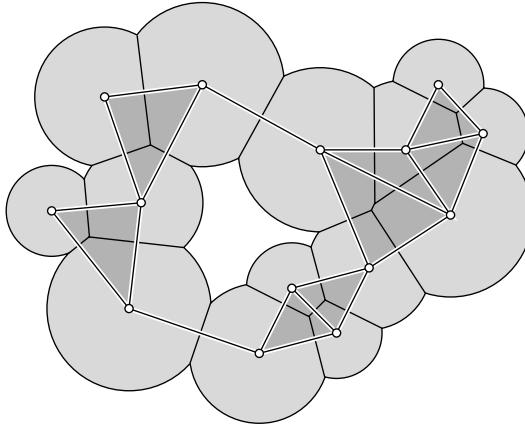


Figure III.16: Convex decomposition of a union of disks and the weighted alpha complex superimposed.

Filtration. There is a free parameter, r , which we may vary to get smaller and larger unions and smaller and larger alpha complexes. Sometimes, there is a best choice of r , but more often there is not. Indeed, the more interesting object is the family of alpha complexes, since it represents the data at different scales, if you will, and it allows us to draw conclusions from comparisons between different complexes in the same family.

We first explain the construction in the relatively straightforward unweighted case. Given a finite set $S \subseteq \mathbb{R}^d$, we continuously increase the radius and thus get a 1-parameter family of nested unions. Correspondingly, we get a 1-parameter family of nested alpha complexes, but because they are all subcomplexes of the same Delaunay complex, which is finite, only finitely many of them are distinct. Writing K_i for the i -th alpha complex in this sequence, we get

$$\emptyset = K_0 \subseteq K_1 \subseteq \dots \subseteq K_m,$$

which we call a *filtration* of $K_m = \text{Delaunay}$. What we have here is a stepwise assembly of the final complex in such a way that every set along the way is a simplicial complex.

There is more than one way to generalize this construction to the weighted case. For example, we could grow the corresponding balls uniformly. Starting with B_u , which has radius $\sqrt{w_u}$, we would increase the radius to $\sqrt{w_u} + r$ for $r > 0$. This makes sense in many applications, including the modeling of biomolecules, but has the complicating side effect that the Voronoi diagram of the set of balls for different values of r are not necessarily the same. Hence, the resulting alpha complexes are not necessarily nested. Instead, we let $B_u(r)$ be the ball with center u and squared radius $w_u + r^2$. The points x with equal power from $B_u(r)$ and $B_v(r)$ satisfy $\|x - u\|^2 - (w_u + r^2) = \|x - v\|^2 - (w_v + r^2)$. The squared radius cancels, implying that the same points x form the bisector for all choices of r . Hence, the union of

balls are decomposed into convex sets by the same weighted Voronoi diagram for any r . Similarly, the weighted alpha complexes are all subcomplexes of the weighted Delaunay triangulation of the given points. More specifically, the alpha complex for r_0 is a subcomplex of that for r whenever $r_0 \leq r$, and we again get a filtration that starts with the empty complex and ends with the entire weighted Delaunay complex, as in the unweighted case.

The structure of a simplex. We are interested in the difference between two contiguous complexes in the filtration, $K_{i+1} - K_i$. For this purpose, we study the structure of an abstract simplex, and not just because it arises as an element of the alpha complex. Recall that an abstract d -simplex, α , is a set of $d + 1$ points. It has 2^{d+1} subsets, including the empty set and α itself. In the *Hasse diagram* of this set system, we draw a node for each subset of α and an arc for each subset relation, avoiding arcs that are implied by transitivity. Drawing the containing sets above the contained ones and keeping subsets of same cardinality in common rows, we get a picture like the one in Figure III.17. It looks like the edge-skeleton of the

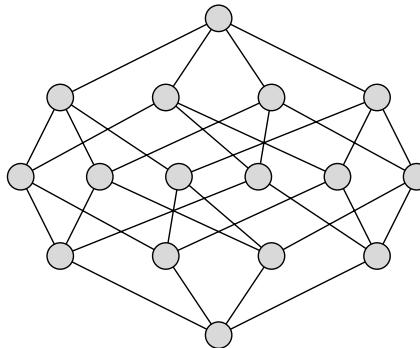


Figure III.17: The Hasse diagram of an abstract 3-simplex. Each containment between faces is represented by a decreasing path connecting the nodes.

$(d + 1)$ -dimensional cube, and not by coincidence. Indeed, we can construct the Hasse diagram inductively, first drawing the Hasse diagram of a $(d - 1)$ -face. By inductive assumption, this is the edge-skeleton of a d -cube. When we add the last point, u_d , to the simplex, we get a new set $\beta \cup \{u_d\}$ for each old set β . To update the Hasse diagram, we add a second copy of the d -cube and connect corresponding sets. This is precisely the recipe for drawing the $(d + 1)$ -cube.

Another useful method for constructing the Hasse diagram is to add one pair of adjacent nodes at a time. We describe this in the other direction, disassembling the diagram one pair at a time. Specifically, we allowed ourselves to remove a pair $\beta_0 \subset \beta$ if β is the only remaining set that properly contains β_0 . Note that β is necessarily maximal and we have $\dim \beta_0 = \dim \beta - 1$ because the operation maintains the system as an abstract simplicial complex. It is easy to see that disassembling the Hasse diagram of the d -simplex this way is possible, for example by removing the pairs $\beta_0 \subset \beta_0 \cup \{u_d\}$ in the order of decreasing dimension of β_0 .

Collapses. Now suppose we have a geometric d -simplex, σ , and we consider the Hasse diagram of its system of faces, to which we add the empty set to be consistent with earlier assumptions. The operation of removing a pair $\beta_0 \subset \beta$ corresponds to removing a pair of faces $\tau_0 < \tau$. The condition is that τ is the only remaining proper coface of τ_0 . The operation of removing the pair $\tau_0 < \tau$ is referred to as an *elementary collapse*, or a $(k, k+1)$ -collapse when $k = \dim \tau_0$. As illustrated in Figure III.18, a d -simplex can be reduced to a single simplex by a sequence of $2^d - 1$ elementary collapses. Since the elementary collapses maintain the set as a simplicial complex, the remaining simplex is necessarily a vertex. We can apply elementary collapses more generally to any simplicial complex, K . Letting L be the result of the collapse, we note that there is a deformation retraction from $|K|$ to $|L|$. This implies that K and L have the same homotopy type. We call K *collapsible* if there is a sequence of elementary collapses that reduces K to a single vertex. Since collapses preserve the homotopy type, this is only possible if $|K|$ is contractible. However, it turns out that not every simplicial complex with contractible underlying space is collapsible.

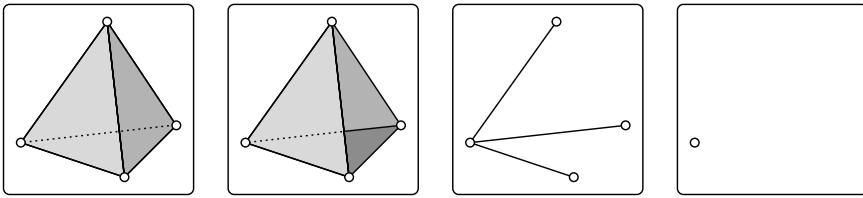


Figure III.18: From left to right: a tetrahedron, the three triangles left after a $(2,3)$ -collapse, the three edges left after three additional $(1,2)$ -collapses, and the vertex left after three additional $(0,1)$ -collapses.

It is convenient to extend the notion of collapse and consider pairs of simplices $\tau < v$ whose dimensions differ by one or more. Instead of requiring that v be the only proper coface of τ , we now require that all cofaces of τ be faces of v . Letting $k = \dim \tau$ and $\ell = \dim v$, we get $\binom{\ell-k}{i}$ simplices of dimension $i+k$ and therefore a total of $2^{\ell-k} = \sum_{i=0}^{\ell-k} \binom{\ell-k}{i}$ simplices between τ and v , including the two. The Hasse diagram of this set of faces has the structure of an $(\ell-k-1)$ -simplex, which we have seen can be collapsed down to a vertex by a sequence of $2^{\ell-k-1} - 1$ elementary collapses. Each $(i, i+1)$ -collapse in this sequence corresponds to an $(i+k+1, i+k+2)$ -collapse in the sequence that removes the cofaces of τ . We append a $(k, k+1)$ -collapse which finally removes τ together with the last remaining proper coface. We refer to this sequence of $2^{\ell-k-1}$ elementary collapses as a (k, ℓ) -collapse. Since elementary collapses preserve the homotopy type, so do the more general collapses.

Critical and regular events. Let r_i be the smallest radius such that $K_i = \text{Alpha}(r_i)$. A simplex τ belongs to K_{i+1} but not to K_i if the balls with radius r_{i+1} have a non-empty common intersection with the corresponding intersection

of Voronoi cells but the balls with radius r_i do not; see Figure III.19. Assuming general position and $\dim \tau = k$, the intersection of Voronoi cells, $V_\tau = \bigcap_{u \in \tau} V_u$, is a convex polyhedron of dimension $d - k$. By definition of r_{i+1} , the balls $B_u(r_{i+1})$ intersect V_τ in a single point, x .

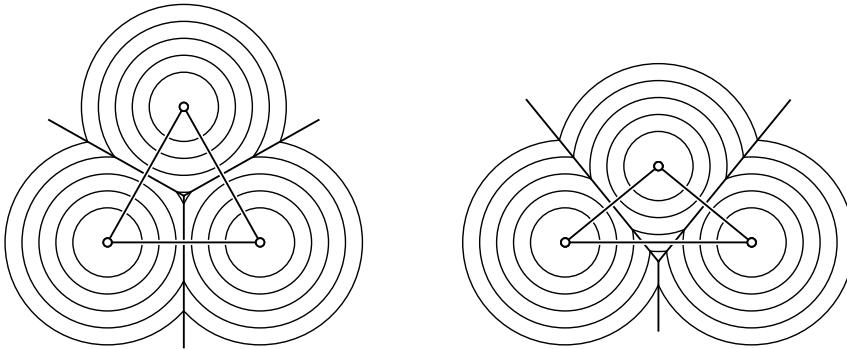


Figure III.19: Left: three points spanning an acute triangle. In the alpha complex evolution, the three edges appear before a critical event adds the triangle. Right: three points spanning an obtuse triangle. Two edges appear before a regular event adds the triangle together with the third edge.

Consider first the case that x lies on the boundary of V_τ . Then there are other Voronoi polyhedra for which x is the first contact with the union of balls. Assume V_τ is the polyhedron with highest dimension in this collection and let V_v be the polyhedron with lowest dimension. Correspondingly, τ is the simplex with lowest dimension in $K_{i+1} - K_i$ and v is the simplex with highest dimension. The other simplices in $K_{i+1} - K_i$ are the faces of v that are cofaces of τ . In other words, we obtain K_i from K_{i+1} by a (k, ℓ) -collapse, where $k = \dim \tau$ and $\ell = \dim v$. We call this collapse a *regular event* in the evolution of the alpha complex.

Consider second the case that x lies in the interior of V_τ and it is not the first contact for any higher-dimensional Voronoi polyhedron. In other words, τ is the only simplex in $K_{i+1} - K_i$. We call the addition of τ a *critical event* because it changes the homotopy type of the complex. Since the union of balls has the homotopy type of the complex, we know that the union also changes its type when the radius reaches r_{i+1} .

Bibliographic notes. Alpha complexes were introduced for points in \mathbb{R}^2 by Edelsbrunner, Kirkpatrick, and Seidel [59], extended to \mathbb{R}^3 in [63], and to weighted points in general, fixed dimension in [53]. The 3-dimensional software written by Ernst Mücke has been popular in many areas of science and engineering, including structural molecular biology where alpha complexes serve as an efficient representation of proteins and other biomolecules. Alpha complexes were the starting point of the work on persistent homology, to be discussed in Chapter VII. The difference between critical and regular events in the evolution of the alpha complex reminds

us of the difference between critical and regular points of a Morse function, which will be studied in Chapter VI. The connection is direct but is made technically difficult because Morse theory was developed principally for smooth functions [111]. A lesser known development of the same ideas for non-smooth functions is based on the concept of a topological Morse function [114] of which the Euclidean distance and power functions for a finite point set are examples.

Exercises

The credit assignment reflects a subjective assessment of difficulty. A typical question can be answered using knowledge of the material combined with some thought and analysis.

1. **Deciding isomorphism** (three credits). What is the computational complexity of recognizing isomorphic abstract simplicial complexes?
2. **Order complex** (two credits). A *flag* in a simplicial complex K in \mathbb{R}^d is a nested sequence of proper faces, $\sigma_0 < \sigma_1 < \dots < \sigma_k$. The collection of flags forms an abstract simplicial complex A sometimes referred to as the *order complex* of K . Prove that A has a geometric realization in \mathbb{R}^d .
3. **Barycentric subdivision** (one credit). Let K consist of a d -simplex σ and its faces.
 - (i) How many d -simplexes belong to the barycentric subdivision, $\text{Sd } K$?
 - (ii) What is the d -dimensional volume of the individual d -simplices in $\text{Sd } K$?
4. **Covering a tree** (one credit). Let P be a finite collection of closed paths that cover a tree; that is, each node and each edge of the tree belongs to at least one path.
 - (i) Prove that the nerve of P is contractible.
 - (ii) Is the nerve still contractible if we allow subtrees in the collection? What about subforests?
5. **Nerve of stars** (one credit). Let K be a simplicial complex. Prove that K is a geometric realization of the nerve of the collection of vertex stars in K .
6. **Helly for boxes** (two credits). A *box* in \mathbb{R}^d is defined by pairs $a_i \leq b_i$ and consists of all points $x = (x_1, x_2, \dots, x_d)$ satisfying $a_i \leq x_i \leq b_i$ for $1 \leq i \leq d$. Let F be a finite collection of boxes in \mathbb{R}^d . Prove that if every pair of boxes has a non-empty intersection, then the entire collection has a non-empty intersection.
7. **Alpha complexes** (two credits). Let $S \subseteq \mathbb{R}^d$ be a finite set of points in general position. Recall that $\check{\text{C}}\text{ech}(r)$ and $\text{Alpha}(r)$ are the $\check{\text{C}}\text{ech}$ and alpha complexes for radius $r \geq 0$. Is it true that $\text{Alpha}(r) = \check{\text{C}}\text{ech}(r) \cap \text{Delaunay}$? If yes, prove the following two subcomplex relations. If no, give examples to show which subcomplex relations are not valid.

- (i) $\text{Alpha}(r) \subseteq \check{\text{Cech}}(r) \cap \text{Delaunay}$.
 - (ii) $\check{\text{Cech}}(r) \cap \text{Delaunay} \subseteq \text{Alpha}(r)$.
8. **Collapsibility** (two credits). Call a simplicial complex *collapsible* if there is a sequence of collapses that reduces the complex to a single vertex. The existence of such a sequence implies that the underlying space of the complex is contractible. Describe a finite 2-dimensional simplicial complex that is not collapsible although its underlying space is contractible.

Part B

Computational Algebraic Topology

Chapter IV

Homology

Homology is a mathematical formalism for talking in a quantitative and unambiguous manner about how a space is connected. Compared to most other, competing formalisms, homology has faster algorithms but captures less of the topological information. We should keep in mind, however, that detailed classifications are not within our computational reach in any case. Specifically, the question of whether or not two triangulated 4-manifolds are homeomorphic or homotopy equivalent are both undecidable. In practice, having fast algorithms is a definitive advantage and being insensitive to some topological information is not necessarily a drawback. More useful than knowing everything is being able to assess the importance of information and to rank it accordingly, a topic we will address directly in Chapter VII. Before we get there, we need to learn the basics, which we do in this chapter.

IV.1 Homology Groups

Homology groups provide a mathematical language for the holes in a topological space. Perhaps surprisingly, they capture holes indirectly, by focusing on what surrounds them. Their main ingredients are group operations and maps that relate topologically meaningful subsets of a space to each other. In this section, we introduce the various groups involved in the setup of homology.

Chain complexes. Let K be a simplicial complex and p a dimension. A p -chain is a formal sum of p -simplices in K . The standard notation for this is $c = \sum a_i \sigma_i$, where the σ_i are the p -simplices and the a_i are the *coefficients*. In computational topology, we mostly work with coefficients a_i that are either 0 or 1, called *modulo 2 coefficients*. Coefficients can, however, be more complicated numbers like integers, rational numbers, real numbers, elements of a field, or elements of a ring. Since we work modulo 2, we can think of a chain as a set of p -simplices, namely those σ_i with $a_i = 1$. But when we do consider chains with other coefficient groups, this way of thinking is more cumbersome, so we will use it sparingly.

Two p -chains are added componentwise, like polynomials. Specifically, if $c = \sum a_i \sigma_i$ and $c' = \sum b_i \sigma_i$, then $c + c' = \sum (a_i + b_i) \sigma_i$, where the coefficients satisfy $1+1=0$. In set notation, the sum of two p -chains is their symmetric difference. The p -chains together with the addition operation form the *group of p -chains* denoted as $(C_p, +)$, or simply $C_p = C_p(K)$ since the operation is understood. Associativity follows from associativity of addition. The neutral element is $0 = \sum 0\sigma_i$. The inverse of c is $-c = c$ since $c + c = 0$. Finally, C_p is abelian because addition modulo 2 is abelian. We have a group of p -chains for each integer p . For p less than zero and greater than the dimension of K this group is trivial, consisting only of the neutral element. To relate these groups, we define the *boundary* of a p -simplex as the sum of its $(p-1)$ -dimensional faces. Writing $\sigma = [u_0, u_1, \dots, u_p]$ for the simplex spanned by the listed vertices, its boundary is

$$\partial_p \sigma = \sum_{j=0}^p [u_0, \dots, \hat{u}_j, \dots, u_p],$$

where the hat indicates that u_j is omitted. For a p -chain, $c = \sum a_i \sigma_i$, the boundary is the sum of the boundaries of its simplices, $\partial_p c = \sum a_i \partial_p \sigma_i$. Hence, taking the boundary maps a p -chain to a $(p-1)$ -chain, and we write $\partial_p : C_p \rightarrow C_{p-1}$. Notice also that taking the boundary commutes with addition, that is, $\partial_p(c + c') = \partial_p c + \partial_p c'$. This is the defining property of a *homomorphism*, which is a map between groups that commutes with the group operation. Therefore, we refer to ∂_p as the *boundary homomorphism* or, for short, the *boundary map* for chains. The *chain complex* is the sequence of chain groups connected by boundary homomorphisms,

$$\dots \xrightarrow{\partial_{p+2}} C_{p+1} \xrightarrow{\partial_{p+1}} C_p \xrightarrow{\partial_p} C_{p-1} \xrightarrow{\partial_{p-1}} \dots.$$

It will often be convenient to drop the index from the boundary homomorphism since it is implied by the dimension of the chain it applies to.

Cycles and boundaries. We distinguish two particular types of chains and use them to define homology groups. A p -cycle is a p -chain with empty boundary, $\partial c = 0$. Since ∂ commutes with addition, we have a *group of p -cycles*, denoted as $Z_p = Z_p(K)$, which is a subgroup of the group of p -chains. In other words, the group of p -cycles is the kernel of the p -th boundary homomorphism, $Z_p = \ker \partial_p$. Since the chain groups are abelian, so are their cycle subgroups. Consider $p = 0$ as an example. The boundary of every vertex is zero, $C_{-1} = 0$; hence $Z_0 = \ker \partial_0 = C_0$. For $p > 0$, however, Z_p is usually not all of C_p .

A p -boundary is a p -chain that is the boundary of a $(p+1)$ -chain, $c = \partial d$ with $d \in C_{p+1}$. Since ∂ commutes with addition, we have a *group of p -boundaries*, denoted by $B_p = B_p(K)$, which is again a subgroup of the p -chains. In other words, the group of p -boundaries is the image of the $(p+1)$ -st boundary homomorphism, $B_p = \text{im } \partial_{p+1}$. Since the chain groups are abelian, so are their boundary subgroups. Consider $p = 0$ as an example. Every 1-chain consists of some number of edges, each with two endpoints. Taking the boundary cancels duplicate endpoints in pairs, leaving an even number of distinct vertices. Now suppose the complex is connected.

Then for any even number of vertices, we can find paths that connect them pairwise and we can add the paths to get a 1-chain whose boundary consists of the given vertices. Hence, every even set of vertices is a 0-boundary and every odd set of vertices is not. If K is connected, this implies that exactly half the 0-cycles are 0-boundaries. The fundamental property that makes homology work is that the boundary of a boundary is necessarily zero.

FUNDAMENTAL LEMMA OF HOMOLOGY. $\partial_p \partial_{p+1} d = 0$ for every integer p and every $(p+1)$ -chain d .

PROOF. We just need to show that $\partial_p \partial_{p+1} \tau = 0$ for a $(p+1)$ -simplex τ . The boundary, $\partial_{p+1} \tau$, consists of all p -faces of τ . Every $(p-1)$ -face of τ belongs to exactly two p -faces, so $\partial_p(\partial_{p+1} \tau) = 0$. \square

It follows that every p -boundary is also a p -cycle or, equivalently, that B_p is a subgroup of Z_p . Figure IV.1 illustrates the subgroup relations among the three types of groups and their connection across dimensions established by the boundary homomorphisms.

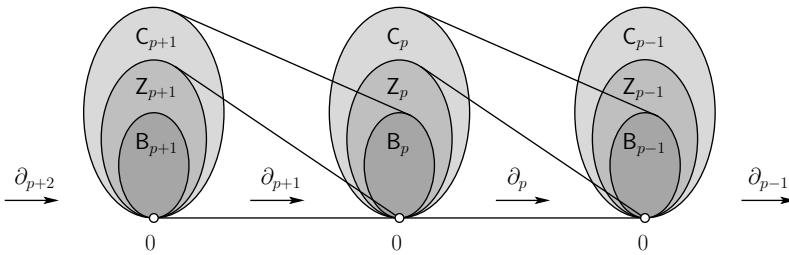


Figure IV.1: The chain complex consisting of a linear sequence of chain, cycle, and boundary groups connected by homomorphisms.

Homology groups. Since the boundaries form subgroups of the cycle groups, we can take quotients. In other words, we can partition each cycle group into classes of cycles that differ from each other by boundaries. This leads to the notion of homology groups and their ranks, which we now define and discuss.

DEFINITION. The p -th homology group is the p -th cycle group modulo the p -th boundary group, $H_p = Z_p / B_p$. The p -th Betti number is the rank of this group, $\beta_p = \text{rank } H_p$.

Each element of $H_p = H_p(K)$ is obtained by adding all p -boundaries to a given p -cycle, $c + B_p$ with $c \in Z_p$. In group theory, $c + B_p$ is called a *coset* of B_p in Z_p . If we take any other cycle $c' = c + c''$, with c'' an element of B_p , we get the same class, $c' + B_p = c + B_p$, since $c'' + B_p = B_p$. This class is thus a coset of H_p and is referred to as a *homology class*. Any two cycles in the same homology class are said

to be *homologous*, which is denoted as $c \sim c'$. We may take c as the representative of this class, but we could choose any other cycle in the class as well. Similarly, addition of two classes, $(c + B_p) + (c_0 + B_p) = (c + c_0) + B_p$, is independent of the representatives and is therefore well defined. We thus see that H_p is indeed a group, and because Z_p is abelian, so is H_p .

The cardinality of a group is called its *order*. Since we use modulo 2 coefficients, a group with n generators has order 2^n . For example, the base 2 logarithm of the order of C_p is the number of p -dimensional simplices in the complex. Furthermore, the group is isomorphic to \mathbb{Z}_2^n , the group of bit-vectors of length n together with the exclusive-or operation. This is an n -dimensional vector space generated by n bit-vectors, for example the n unit vectors. The dimension is referred to as the *rank* of the vector space, $n = \text{rank } \mathbb{Z}_2^n = \log_2 \text{ord } \mathbb{Z}_2^n$. The cycles and boundaries exhibit the same vector space structure, except that their dimension is often less than that of the chains. The number of cycles in each homology class is the order of B_p ; hence the number of classes in the homology group is $\text{ord } H_p = \text{ord } Z_p / \text{ord } B_p$. Equivalently, the rank is the difference, $\beta_p = \text{rank } H_p = \text{rank } Z_p - \text{rank } B_p$. This

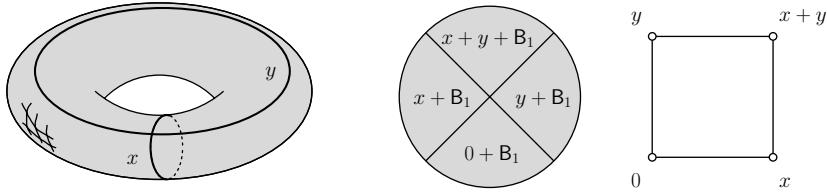


Figure IV.2: The first homology group of the torus has order 4 and rank 2. In the middle, the four elements are drawn as the cosets in the group of 1-cycles. On the right, the four elements are the vertices of a square.

suggests two alternative methods for illustrating a homology group, as a partition of the set of cycles or the hypercube of dimension β_p . As an example consider the torus in Figure IV.2. There are only four homology classes in H_1 , namely B_1 , $x + B_1$, $y + B_1$, and $(x + y) + B_1$, where x and y are the non-bounding 1-cycles that go once around the arm and the hole of the torus. The two corresponding cosets, $x + B_1$ and $y + B_1$, generate the first homology group.

The homology of a ball. We define a closed ball to be any triangulated topological space that is homeomorphic to \mathbb{B}^k , the subset of points at a distance at most one from the origin in \mathbb{R}^k . What is the homology of a ball? Given our intuition that homology should measure holes, it should be trivial. This almost turns out to be true; actually if K triangulates a ball, then $H_p(K) = 0$ except when $p = 0$ where it has rank 1. This is surprisingly hard to prove, however! It is usually done with a lot of machinery like simplicial approximations and homotopy equivalences. For now, let us at least see this directly when K is the set of faces of a single simplex of dimension k . In this case, the p -chains of K have rank equal to the number of p -faces, which is $\binom{k+1}{p+1}$. Let the vertices be u_0, u_1, \dots, u_k and consider a p -chain c with simplices of the form $[u_{i_0}, u_{i_1}, \dots, u_{i_p}]$. The condition $\partial c = 0$ is equivalent to

every $(p-1)$ -simplex occurring an even number of times as a face of p -simplices in c . Assuming $p > 0$, we can construct a $(p+1)$ -chain d with boundary $\partial d = c$. This will imply that every p -cycle is also a p -boundary and, equivalently, that H_p is trivial. Specifically, we let d be the set of $(p+1)$ -simplices of the form $[u_0, u_{i_0}, u_{i_1}, \dots, u_{i_p}]$. In words, d contains a $(p+1)$ -simplex for each p -simplex in c that does not have u_0 as a vertex. To see that c is the boundary of d , we distinguish three types of p -faces of simplices in d . A p -simplex in c that does not contain u_0 occurs exactly once as a face of a $(p+1)$ -simplex in d and therefore belongs to ∂d . A p -simplex τ not in c occurs an even number of times, namely once for each time the $(p-1)$ -face σ obtained by dropping u_0 occurs in the boundary of a simplex in c . By the same argument, we get a p -simplex τ in c that contains u_0 an odd number of times because one of the p -simplices in c that contains the $(p-1)$ -face σ does not give rise to a $(p+1)$ -simplex in d , namely τ itself.

This covers all positive dimensions. For $p = 0$, we have already observed that exactly half the cycles are boundaries. Hence, $H_0 = Z_0/B_0$ is isomorphic to \mathbb{Z}_2 and $\beta_0 = 1$, as claimed.

Reduced homology. There is something dissatisfying about the 0-th homology group behaving differently for the ball than the others. The reason for the difference is that we have set up things so that β_0 counts the components, but if there is one component, there is no hole. More satisfying would be to count one when we have two components, namely for the one gap between them. This is achieved by a small but often useful modification of homology, namely adding the *augmentation map* $\epsilon : C_0 \rightarrow \mathbb{Z}_2$ defined by $\epsilon(u) = 1$ for each vertex u . We thus get

$$\dots \xrightarrow{\partial_2} C_1 \xrightarrow{\partial_1} C_0 \xrightarrow{\epsilon} \mathbb{Z}_2 \xrightarrow{0} 0 \longrightarrow \dots .$$

Cycles and boundaries are defined as before, and the only difference we notice is for Z_0 which now requires that each 0-cycle have an even number of vertices. This results in the *reduced homology groups*, \tilde{H}_p , and the *reduced Betti numbers*, $\tilde{\beta}_p = \text{rank } \tilde{H}_p$. Assuming K is non-empty, we have $\tilde{\beta}_p = \beta_p$ for all $p \geq 1$ and $\tilde{\beta}_0 = \beta_0 - 1$. For $K = \emptyset$, we have $\tilde{\beta}_{-1} = 1$. This is because both elements of \mathbb{Z}_2 belong to the kernel and are therefore (-1) -cycles, but only one belongs to the image of the augmentation map and is therefore a (-1) -boundary.

Induced maps. A continuous map from one topological space to another maps cycles to cycles and boundaries to boundaries. We can therefore use the images to construct new homology groups. These are not necessarily the same as the ones of the original space since cycles can become boundaries, for example trivial cycles. We describe this more formally for two simplicial complexes and a simplicial map, $f : K \rightarrow L$, between them. Recall that f takes each simplex of K linearly to a simplex of L . It induces a map from the chains of K to the chains of the same dimension of L . Specifically, if $c = \sum a_i \sigma_i$ is a p -chain in K , then $f_\#(c) = \sum a_i \tau_i$, where $\tau_i = f(\sigma_i)$ if it has dimension p and $\tau_i = 0$ if $f(\sigma)$ has dimension less than p . Writing ∂_K and ∂_L for the boundary maps in the two complexes, we note that

$f_\# \circ \partial_K = \partial_L \circ f_\#$, that is, the induced map commutes with the boundary maps. This is obvious when $f(\sigma_i)$ has dimension p , since then all $(p-1)$ -faces of σ_i map to the corresponding $(p-1)$ -faces of τ_i . If, on the other hand, $f(\sigma_i)$ has dimension less than p , then the $(p-1)$ -faces of σ_i map to simplices of dimension less than $p-1$, with the possible exception of exactly two $(p-1)$ -faces whose images coincide and cancel each other. So both $f_\#(\partial_K \sigma_i)$ and $\partial_L f_\#(\sigma_i)$ are zero. Note that in the case when $f : K \rightarrow L$ is the inclusion of one simplicial complex into another, simplices always keep their dimension, so the induced map, $f_\#$, is a little easier to understand.

The fact that the induced map commutes with the boundary map implies that $f_\#$ takes cycles to cycles, $f_\#(\mathbb{Z}_p(K)) \subseteq f_\#(\mathbb{Z}_p(L))$, and boundaries to boundaries, $f_\#(\mathbb{B}_p(K)) \subseteq f_\#(\mathbb{B}_p(L))$. Therefore, it defines a map on the quotients, which we call the *induced map on homology*, sometimes written $f_* : \mathbb{H}_p(K) \rightarrow \mathbb{H}_p(L)$. Note that the rank of the image is bounded from above by both Betti numbers, $\text{rank } f_*(\mathbb{H}_p(K)) \leq \min\{\beta_p(K), \beta_p(L)\}$.

Degree of a map. We now present a first application of the concept of induced maps. We describe it for general continuous maps, appealing to the Simplicial Approximation Theorem proved in Section III.1 when we need triangulations and an approximating simplicial map. Let $g : \mathbb{S}^p \rightarrow \mathbb{S}^p$ be a continuous map and let c be the unique generator of the p -th homology group of the p -sphere. Then $g(c)$ is either homologous to c or to 0. In other words, $g(c) \sim \alpha c$ and $\alpha \in \{0, 1\}$ is called the *modulo 2 degree* of g . If g is the identity, then $\alpha = 1$. However, if g extends a continuous map $g_0 : \mathbb{B}^{p+1} \rightarrow \mathbb{S}^p$, then the induced map on homology, $g_* : \mathbb{H}_p(\mathbb{S}^p) \rightarrow \mathbb{H}_p(\mathbb{S}^p)$, is the composite of two induced maps, $\mathbb{H}_p(\mathbb{S}^p) \rightarrow \mathbb{H}_p(\mathbb{B}^{p+1}) \rightarrow \mathbb{H}_p(\mathbb{S}^p)$, where the first is induced by inclusion. The middle group is trivial; hence $\alpha = 0$. We are now ready to prove a classic result on fixed points of continuous maps.

BROUWER'S FIXED POINT THEOREM. A continuous map $f : \mathbb{B}^{p+1} \rightarrow \mathbb{B}^{p+1}$ has at least one fixed point $x = f(x)$.

PROOF. Let $A, B : \mathbb{S}^p \rightarrow \mathbb{S}^p$ be maps defined by $A(x) = (x - f(x))/\|x - f(x)\|$ and $B(x) = x$. Since B is the identity, its modulo 2 degree is 1. If f has no fixed point, then A is well defined and has modulo 2 degree 0 because it extends a map from the $(p+1)$ -ball to the p -sphere. We now construct $H : \mathbb{S}^p \times [0, 1] \rightarrow \mathbb{S}^p$ by setting $H(x, t) = (x - tf(x))/\|x - tf(x)\|$. For $t = 1$, we have $x \neq f(x)$ because there is no fixed point, and for $t < 1$, we have $x \neq tf(x)$ because $\|x\| = 1 > \|tf(x)\|$. We conclude that H is a homotopy between A and B , which implies that the modulo 2 degree of the two are the same, a contradiction. \square

Bibliographic notes. Like many other concepts in topology, homology groups were introduced by Henri Poincaré in one of a series of papers on “analysis situ” [121]. He named the ranks of the homology groups after another mathematician, Betti, who introduced a slightly different version years earlier. The field experienced a rapid development during the twentieth century. There were many competing

theories, simplicial and singular homology being just two examples, which have been consolidated by axiomatizing the assumptions under which homology groups exist [66]. Today we have a number of well-established textbooks in the field. We refer to Giblin [78] for an intuitive introduction and to Hatcher [82], Munkres [116], and Spanier [136] for more comprehensive sources. Brouwer's Fixed Point Theorem impresses by its generality and is popular outside of mathematics also. He proved the 3-dimensional case in 1910 [22] and the general case in 1912 [23].

IV.2 Matrix Reduction

The homology groups of a triangulated space can be computed from the matrices representing the boundary homomorphisms. Their reduced versions readily provide the ranks of the cycle and boundary groups, and their differences give the Betti numbers. Summing these same differences leads to a proof of the Euler-Poincaré formula which generalizes the Euler relation for planar graphs.

Euler-Poincaré formula. Recall that the Euler characteristic of a simplicial complex is the alternating sum of the number of simplices in each dimension. Also recall that the rank of the p -th homology group is the rank of the p -th cycle group minus the rank of the p -th boundary group. Writing $z_p = \text{rank } Z_p$ and $b_p = \text{rank } B_p$, this can be stated as $\beta_p = z_p - b_p$. Furthermore, writing $n_p = \text{rank } C_p$ for the number of p -simplices in K , we know that $n_p = z_p + b_{p-1}$. This is the general fact that for any linear transformation between vector spaces $f : U \rightarrow V$, the dimension of U equals the sum of the dimension of the kernel of f and the dimension of the image of f . The Euler characteristic is the alternating sum of the n_p , which is therefore

$$\begin{aligned}\chi &= \sum_{p \geq 0} (-1)^p(z_p + b_{p-1}) \\ &= \sum_{p \geq 0} (-1)^p(z_p - b_p) \\ &= \sum_{p \geq 0} (-1)^p\beta_p.\end{aligned}$$

To appreciate the beauty of this result, we need to know that homology groups do not depend on the triangulation chosen for a topological space. The technical proof of this claim is difficult, and we refer the reader to more advanced texts. Even the more general result that homotopy equivalent spaces have isomorphic homology groups is plausible. For example, we can free ourselves from the triangulation entirely and define chains in terms of continuous maps from the standard simplex into the space \mathbb{X} . This gives rise to singular homology, which can be shown to give groups isomorphic to the ones we get by simplicial homology, the theory we describe in this chapter. If we now have a continuous map $f : \mathbb{X} \rightarrow \mathbb{Y}$, we can map the chains from \mathbb{X} to those of \mathbb{Y} by simply composing. If f is a homotopy equivalence, then it turns out that \mathbb{X} and \mathbb{Y} have isomorphic homology groups. This also implies that

the Euler characteristic is an invariant of the space, that is, it does not depend on the simplicial complex we use to triangulate it.

EULER-POINCARÉ THEOREM. The Euler characteristic of a topological space is the alternating sum of its Betti numbers, $\chi = \sum_{p \geq 0} (-1)^p \beta_p$.

Boundary matrices. To compute homology, we combine information from two sources, one representing the cycles and the other the boundaries, just as in the proof of the Euler-Poincaré Theorem. Let K be a simplicial complex. Its p -th boundary matrix represents the $(p - 1)$ -simplices as rows and the p -simplices as columns. Assuming an arbitrary but fixed ordering of the simplices, for each dimension, this matrix is $\partial_p = [a_i^j]$, where i ranges from 1 to n_{p-1} , j ranges from 1 to n_p , and $a_i^j = 1$ if the i -th $(p - 1)$ -simplex is a face of the j -th p -simplex and $a_i^j = 0$, otherwise. Given a p -chain, $c = \sum a_i \sigma_i$, the boundary, $\partial_p c$, can be computed by matrix multiplication:

$$\begin{bmatrix} a_1^1 & a_1^2 & \dots & a_1^{n_p} \\ a_2^1 & a_2^2 & \dots & a_2^{n_p} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n_{p-1}}^1 & a_{n_{p-1}}^2 & \dots & a_{n_{p-1}}^{n_p} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{n_p} \end{bmatrix}.$$

In words, a collection of columns represents a p -chain and the sum of these columns gives its boundary. Similarly, a collection of rows represents a $(p - 1)$ -chain and the sum of these rows gives its coboundary, a concept that will be defined in the next chapter.

Row and column operations. The rows of the matrix ∂_p form a basis of the $(p - 1)$ -st chain group, C_{p-1} , and the columns form a basis of the p -th chain group, C_p . We use two types of column operations to modify the matrix without changing its rank: exchanging columns k and l and adding column k to column l . Both can be expressed by multiplying with a matrix $V = [v_i^j]$ on the right. To exchange two columns, we have $v_k^l = v_l^k = 1$ and $v_i^j = 1$ for all $i \neq k, l$. All other entries are

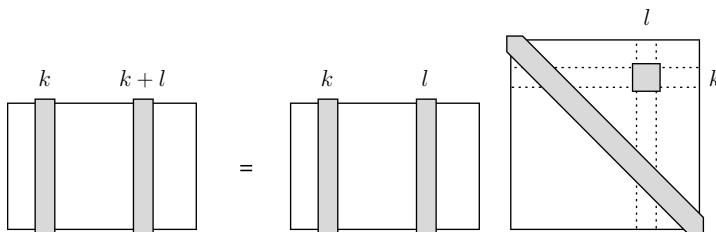


Figure IV.3: The effect of a single off-diagonal one in the matrix V is the addition of one column in the boundary matrix to another. The effect on the basis of C_p is similar.

zero. To add column k to column l , we have $v_k^l = 1$ and $v_i^l = 1$ for all i . All other entries are zero. As indicated in Figure IV.3, the effect of the operation is that the l -th column now represents the sum of the k -th and the l -th p -simplices, or the sum of whatever the two columns represented prior to the operation. Similarly, we have two row operations, one exchanging two rows and the other adding one row to another. This translates to multiplication by a matrix $U = [u_i^j]$ on the left. To exchange two rows, we again have $u_k^l = u_l^k = 1$, $u_i^l = 1$ for $i \neq k, l$, and all other entries zero. To add the k -th to the l -th row, we have $u_l^k = 1$, $u_i^l = 1$ for all i , and all other entries zero, as in Figure IV.4. The effect of this operation is that the k -th

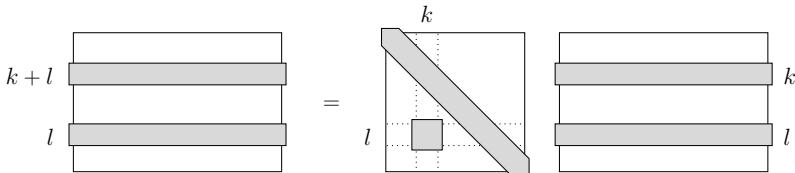


Figure IV.4: The effect of a single off-diagonal one in the matrix U is the addition of one row in the boundary matrix to another. The effect on the basis of C_{p-1} is that the row that was added now represents the sum of $(p-1)$ -chains, the opposite of a column operation.

row now represents the sum of the k -th and the l -th $(p-1)$ -simplices, or the sum of whatever the two rows represented prior to the operation. Although the $(p-1)$ - and p -chains represented by the rows and columns change as we perform row and column operations, they always represent bases of the two chain groups.

Smith normal form. Using row and column operations, we can reduce the p -th boundary matrix to *Smith normal form*. For modulo 2 arithmetic, this means an initial segment of the diagonal is 1 and everything else is 0, as in Figure IV.5. Recall that $n_p = \text{rank } C_p$ is the number of columns of the p -th boundary matrix.

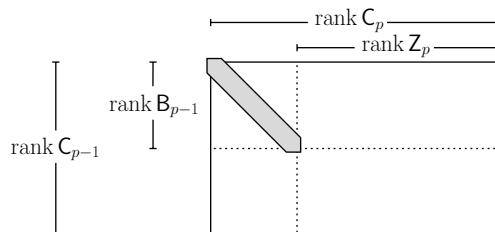


Figure IV.5: The entries in the shaded, initial portion of the diagonal are 1 and all other entries are 0. The ranks of the boundary and cycle groups are readily available as the numbers of non-zero and zero columns.

Let $n_p = b_{p-1} + z_p$ so that the leftmost b_{p-1} columns have ones in the diagonal and the rightmost z_p columns are zero. The former represent p -chains whose non-zero

boundaries generate the group of $(p - 1)$ -boundaries. The latter represent p -cycles that generate Z_p . Once we have all boundary matrices in normal form, we can extract the Betti numbers as differences between ranks, $\beta_p = \text{rank } Z_p - \text{rank } B_p$ for $p \geq 0$. To get the bases of the boundary and cycle groups, we keep track of the matrix products that represent the row and column operations. Writing U_{p-1} and V_p for the left and right products, we get the normal form as $N_p = U_{p-1} \partial_p V_p$. The new basis for the cycle group is given in the last z_p columns of V_p . Similarly, the new basis for the boundary group is encoded in U_{p-1} , and we get the basis vectors from the first b_{p-1} columns of the inverse.

Reduction. To reduce ∂_p , we proceed similarly to Gaussian elimination for solving a system of linear equations. In at most two exchange operations, we move a 1 to the upper left corner, and with at most $n_{p-1} - 1$ row and $n_p - 1$ column additions, we zero out the rest of the first column and first row. We then recurse for the submatrix obtained by removing the first row and first column. We start the reduction by initializing the matrix to $N_p[i, j] = a_i^j$ for all i and j and by calling the function for $x = 1$, the position of the considered diagonal element.

```

void REDUCE( $x$ )
    if there exist  $k \geq x, l \geq x$  with  $N_p[k, l] = 1$  then
        exchange rows  $x$  and  $k$ ; exchange columns  $x$  and  $l$ ;
        for  $i = x + 1$  to  $n_{p-1}$  do
            if  $N_p[i, x] = 1$  then add row  $x$  to row  $i$  endif
        endfor;
        for  $j = x + 1$  to  $n_p$  do
            if  $N_p[x, j] = 1$  then add column  $x$  to column  $j$  endif
        endfor;
        REDUCE( $x + 1$ )
    endif.

```

We have at most n_{p-1} row and n_p column operations per recursive call and hence at most $(n_{p-1} + n_p) \min\{n_{p-1}, n_p\}$ row and column operations in total. Multiplying with their lengths, we thus get a running time of a constant times $2n_{p-1}n_p \min\{n_{p-1}, n_p\}$. The amount of memory is at most some constant times $(n_{p-1} + n_p)^2$ needed to store the matrices. In summary, we reduce the boundary matrices in time at most cubic and in memory at most quadratic in the number of simplices in K .

Example. To get a feeling for the algorithm, we use it to compute the reduced homology group of the 3-ball triangulated by the faces of a single tetrahedron. We do the computations one dimension at a time and this way get the reduced Betti numbers of all skeleta of the complex as we go. The 0-skeleton consists of four vertices, and its sole non-trivial boundary matrix is ∂_0 , consisting of a row of ones, shown as part of the first equation in Figure IV.6. Three column operations remove three of the four ones, and we get $\tilde{\beta}_0 = 3$, the number of zero columns in N_0 . Proceeding to the 1-skeleton, we add the six edges and consider the first boundary

The figure consists of four rows of equations, each showing a boundary matrix (shaded) and a normal form matrix (white) connected by an equals sign. To the right of each equation is a circuit diagram illustrating the reduction process.

- Top row:** Shows the reduction of the zeroth boundary matrix. The boundary matrix has bases a, b, c, d and is shaded. The normal form matrix is white. The circuit diagram shows a simple connection between the two.
- Second row:** Shows the reduction of the first boundary matrix. The boundary matrix has bases ab, ac, ad, bc, bd, cd and is shaded. The normal form matrix has bases a, b, c, d and is white. The circuit diagram shows a more complex connection involving multiple nodes and loops.
- Third row:** Shows the reduction of the second boundary matrix. The boundary matrix has bases abc, abd, acd, bcd and is shaded. The normal form matrix has bases ab, ac, ad, bc, bd, cd and is white. The circuit diagram is very complex, showing many nodes and loops.
- Bottom row:** Shows the reduction of the third boundary matrix. The boundary matrix has bases $abcd$ and is shaded. The normal form matrix has bases abc, abd, acd, bcd and is white. The circuit diagram is similar to the one in the third row.

Figure IV.6: From top to bottom: the matrix equations $N_p = U_{p-1} \partial_p V_p$ for reducing the zeroth, first, second, and third boundary matrices of the tetrahedron. The ones are shaded, and the zeros are white. The bases are indicated for both the boundary and the normal form matrices. For clarity, no exchanges are performed.

matrix. After reduction, it has three ones in the diagonal, shown as part of the second equation in Figure IV.6. Combining the information from N_0 and N_1 , we get $\tilde{\beta}_0 = 3 - 3 = 0$, and counting the zero columns in N_1 , we get $\tilde{\beta}_1 = 3$. Proceeding to the 2-skeleton, we add the four triangles and thus get a triangulation of the 2-sphere. After reduction, the boundary matrix again has three ones in the diagonal, shown as part of the third equation in Figure IV.6. The triangles do not affect the zeroth homology group, and we have $\tilde{\beta}_0 = 0$, as before. Combining the information from N_1 and N_2 , we get $\tilde{\beta}_1 = 3 - 3 = 0$, and counting the zero columns in N_2 , we get $\tilde{\beta}_2 = 1$. We finally get the triangulation of the 3-ball by adding the one tetrahedron. After reduction, the boundary matrix has a single one in the diagonal, shown as part of the fourth equation in Figure IV.6. The first two reduced Betti numbers remain unaffected, and the other two also vanish, so we get $\tilde{\beta}_0 = \tilde{\beta}_1 = \tilde{\beta}_2 = \tilde{\beta}_3 = 0$, as expected.

Bibliographic notes. The generalization of the Euler relation for planar graphs to the Euler-Poincaré Theorem has an interesting history, analyzed from a philosophical viewpoint by Lakatos [99]. The Smith normal form for coefficient groups other than \mathbb{Z}_2 is given in [134]. Already for integers, this complicates matters significantly, and it is no longer straightforward to guarantee a running time that is polynomial in the number of simplices [116]. However, improvements to polynomial time are possible [88, 138].

IV.3 Relative Homology

We extend homology beyond closed spaces by considering nested pairs of closed spaces and studying their difference. We use induced maps on homology to relate the homology of such a pair to the homology of the individual closed spaces.

Relative homology groups. In this chapter, homology groups have been defined for triangulated spaces, which are therefore necessarily closed. To extend them to other spaces, we introduce homology groups for pairs of closed spaces. Let K be a simplicial complex and K_0 a subcomplex of K . The *relative chain groups* are quotients of the chain groups of K and of K_0 ; that is, $C_p(K, K_0) = C_p(K)/C_p(K_0)$. Taking this quotient partitions $C_p(K)$ into cosets, $c + C_p(K_0)$, whose p -chains possibly differ in the p -simplices in K_0 but not in the ones in $K - K_0$. The *boundary map* is induced by the one for K ; that is, $\partial_p(c + C_p(K_0)) = \partial_p c + C_{p-1}(K_0)$. As before, ∂ commutes with addition, and taking the boundary twice gives zero. We thus define *relative cycle groups*, *relative boundary groups*, and *relative homology groups* as kernels, images, and quotients,

$$\begin{aligned} Z_p(K, K_0) &= \ker \partial_p; \\ B_p(K, K_0) &= \text{im } \partial_{p+1}; \\ H_p(K, K_0) &= \ker \partial_p / \text{im } \partial_{p+1}, \end{aligned}$$

just as before. Let $c + C_p(K_0)$ be a relative p -chain. It is a relative p -cycle iff ∂c is carried by K_0 , which includes the possibility that ∂c is zero. Furthermore, it is a relative p -boundary if there is a $(p+1)$ -chain d of K such that $c - \partial d$ is carried by K_0 . As before, we call two relative cycles homologous if their difference is a boundary. To distinguish the homology of a complex K from that of a pair (K, K_0) , we sometimes refer to $H_p(K)$ as the *absolute homology* of K and to its elements as *absolute classes*.

Recall that continuous maps between spaces induce maps on homology. This idea extends to relative homology. Suppose we have a simplicial map, $f : K \rightarrow L$, and subcomplexes $K_0 \subseteq K$ and $L_0 \subseteq L$ such that K_0 maps into L_0 . Then we have an induced map, $f_\#$, between the chains of (K, K_0) and of (L, L_0) . Furthermore, $f_\#$ commutes with ∂ , so it induces a map on homology, $f_* : H_p(K, K_0) \rightarrow H_p(L, L_0)$. A useful example of such an induced map on homology is from (K, \emptyset) to (K, K_0) , as we will see next.

Examples. Relative homology is less intuitive than absolute homology. To help us understand it better, it is useful to compare the two. The difference is described in the algebra of the induced map on homology from $H_p(K) = H_p(K, \emptyset)$ to $H_p(K, K_0)$. Some classes of $H_p(K)$ die entering $H_p(K, K_0)$, some survive, and some new ones get born. We present examples to illustrate how these changes occur.

Let K triangulate the annulus, $\mathbb{S}^1 \times [0, 1]$. In Figure IV.7 on the left, we assume K_0 triangulates the boundary of the annulus, which consists of the circles $\mathbb{S}^1 \times 0$ and $\mathbb{S}^1 \times 1$. First notice that if u is a vertex of K , then either u lies in K_0 or there is a path from u to a vertex u_0 in K_0 . In the second case, the 1-chain traced out by the path has boundary equal to $u + u_0$. It follows that any 0-chain in K is homologous to one in K_0 , which implies that $H_0(K, K_0)$ is trivial. Since $H_0(K)$ has rank one, this is an example of how an absolute homology class dies in relative homology. Next, suppose we have a relative 1-cycle, c . From the definition, we know that its

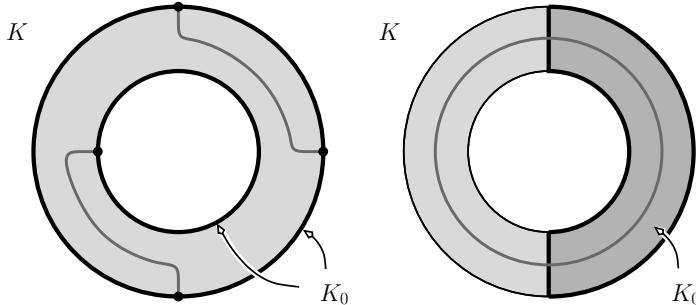


Figure IV.7: The annulus triangulated by a complex K . Left: $K_0 \subseteq K$ triangulates the boundary. The upper right path is a relative boundary and the lower left path generates a non-trivial class in the first relative homology group. Right: $K_0 \subseteq K$ triangulates the right half of the annulus. The core circle generates a non-trivial class in the first absolute homology group as well as the first relative homology group.

boundary, ∂c , is a sum of an even number of vertices in K_0 . In fact, c is a sum of closed loops and of paths that connect two vertices in K_0 . Consider first the case that c is a loop, that is, a 1-cycle in K . Then c is either homologous to zero or to the core circle, $\mathbb{S}^1 \times \frac{1}{2}$. But the core circle is also homologous to $\mathbb{S}^1 \times 0$, which lies in K_0 . Thus, any absolute 1-cycle is trivial in relative homology; that is, the induced map from $H_1(K)$ to $H_1(K, K_0)$ is zero. The core circle therefore gives us another example of an absolute class that dies in relative homology. Consider second the case that c is a path whose endpoints, u_0 and u_1 , lie in K_0 . If both points lie in the same component, either the outer or the inner boundary circle, then there is a 2-chain whose boundary is $c + c_0$, where c_0 is a path in K_0 . On the other hand, if u_0 and u_1 lie in different components, then no such 2-chain exists. Any two such paths are homologous, however, so $H_1(K, K_0)$ has rank one. The path connecting the two components of K_0 is therefore an example of a relative class that is not an absolute class. Finally, $H_2(K)$ is 0, but $H_2(K, K_0)$ has rank one. The new class is the sum of all triangles in K , and its boundary is the sum of all edges in K_0 .

A trivial example of a situation in which a class can be both absolute and relative is given by taking $K_0 = \emptyset$. For a more interesting example, let K_0 triangulate half the annulus, as illustrated in Figure IV.7 on the right. Here, the core circle generates the sole non-trivial class in $H_1(K)$ as well as the one in $H_1(K, K_0)$. In this case, the zeroth and second relative homology groups are both zero and $H_1(K, K_0)$ is the only non-trivial relative homology group of the pair.

Excision. By construction, relative homology depends only on the part of K outside K_0 and ignores the part inside K_0 . Hence, we can remove simplices from both complexes without changing the homology.

EXCISION THEOREM. Let $K_0 \subseteq K$ and $L_0 \subseteq L$ be pairs of simplicial complexes that satisfy $L \subseteq K$ and $L - L_0 = K - K_0$. Then they have isomorphic relative homology groups; that is, $H_p(K, K_0) \simeq H_p(L, L_0)$ for all dimensions p .

Instead of giving a direct algebraic proof of this fact, we take a look at the Smith normal form reduction for relative homology. Ordering the simplices in K_0 before the ones in $K - K_0$, all the relevant information is contained in the lower right submatrices that belong to rows and columns of simplices in $K - K_0$. We reduce these submatrices, ignoring the rows and columns of simplices in K_0 . As illustrated in Figure IV.8, we get the ranks of the relative boundary and cycle groups by counting the non-zero and zero columns in the submatrices. Using the same ordering of simplices, we get the boundary matrices of L by removing rows and columns that correspond to simplices in $K - L$. By definition of L and L_0 , these rows and columns correspond to simplices in K_0 . The lower right submatrices defined by $L - L_0$ are therefore the same as before. This implies $H_p(K, K_0) \simeq H_p(L, L_0)$ for all dimensions p , as claimed in the Excision Theorem.

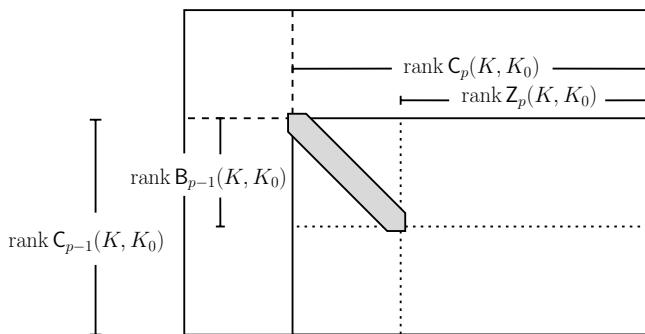


Figure IV.8: By ordering the simplices of K_0 before the others, we get the incidences between simplices in $K - K_0$ in the lower right submatrix, which we reduce to compute the homology of the pair (K, K_0) .

We could have deleted the rows and columns of simplices in K_0 but chose to keep them because they contain the information that relates the relative homology

groups of (K, K_0) with the (absolute) homology groups of K and K_0 . We need new concepts to describe this connection.

Maps between vector spaces. Since we use modulo 2 arithmetic, the induced map on homology is a linear transformation between vector spaces. In general, if $f : U \rightarrow V$ is a linear transformation between vector spaces, the *kernel*, *image*, and *cokernel* are defined as usual:

$$\begin{aligned}\ker f &= \{u \in U \mid f(u) = 0 \in V\}; \\ \text{im } f &= \{v \in V \mid \text{there exists } u \in U \text{ with } f(u) = v\}; \\ \text{cok } f &= V/\text{im } f.\end{aligned}$$

For example, if f is represented by a matrix, such as ∂ , we can reduce to get the kernel spanned by the zero columns, the image spanned by the non-zero rows, and the cokernel spanned by the zero rows. All three are vector spaces in their own right, so we can take direct sums, recalling that this is like taking Cartesian products and using the group operations componentwise. A fundamental result from linear algebra states that U and V are completely described by the three. Specifically, U is isomorphic to the direct sum of the kernel and the image, and V is isomorphic to the direct sum of the image and the cokernel:

$$\begin{aligned}U &\simeq \ker f \oplus \text{im } f; \\ V &\simeq \text{im } f \oplus \text{cok } f.\end{aligned}$$

Again, this has obvious interpretations in terms of the reduced matrix representing f . If we have three vector spaces and two linear transformations, $f : U \rightarrow V$ and $g : V \rightarrow W$, we say the sequence $U \rightarrow V \rightarrow W$ is *exact* at V if $\text{im } f = \ker g$; see Figure IV.9. Note that this implies $g \circ f = 0$; thus the sequence might be three terms in a chain complex, and exactness would mean that the homology group at V was 0. But we will use this concept in more general ways than that. If $0 \rightarrow U \rightarrow V$

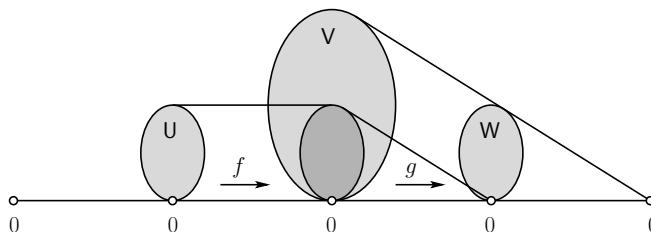


Figure IV.9: A short exact sequence of vector spaces. It starts and ends with zero and is exact at each of the three vector spaces between the two ends.

is a sequence, then exactness at U is equivalent to injectivity of $U \rightarrow V$. Similarly, for $V \rightarrow W \rightarrow 0$, exactness at W is equivalent to surjectivity of $V \rightarrow W$. A *short exact sequence* is a sequence of length five,

$$0 \rightarrow U \xrightarrow{f} V \xrightarrow{g} W \rightarrow 0,$$

that starts and ends with the trivial vector space and is exact at U , V , and W . By what we said above, $f : U \rightarrow V$ is injective and $g : V \rightarrow W$ is surjective. In this situation, it is always true that the middle vector space is isomorphic to the direct sum of the adjacent vector spaces, $V \simeq U \oplus W$. Thus if we somehow already know U and W , then we have calculated V .

Exact sequence of a pair. Sequences that are exact are a convenient means to express otherwise cumbersome relationships between homology groups. Exceptionally powerful are *long exact sequences* which are infinite sequences of vector spaces that are exact at each term. A long exact sequence is like a chain complex, but with trivial homology throughout. A particular example relates the relative homology groups of a pair to the absolute homology groups of the spaces forming the pair.

EXACT SEQUENCE OF A PAIR THEOREM. Let K be a simplicial complex and $K_0 \subseteq K$ a subcomplex. Then there is a long exact sequence

$$\dots \rightarrow H_p(K_0) \rightarrow H_p(K) \rightarrow H_p(K, K_0) \rightarrow H_{p-1}(K_0) \rightarrow \dots .$$

The statement also holds if we substitute the reduced homology groups of K and K_0 for their non-reduced homology groups.

The next section will give a general method for constructing long exact sequences, including that of a pair. Here, we will content ourselves with a brief description of the maps between the groups. The map $H_p(K_0) \rightarrow H_p(K)$ is just the map on homology induced by the inclusion $K_0 \subseteq K$. The map $H_p(K) \rightarrow H_p(K, K_0)$ is also induced by inclusion, $K \subseteq K$; that is, a class generated by a p -cycle c in K is mapped to the relative class generated by $c + C_p(K_0)$. The third map, $H_p(K, K_0) \rightarrow H_{p-1}(K_0)$, is called the *connecting homomorphism* and is the crucial piece of the construction. To describe it, let $c = \sum a_i \sigma_i$ generate a relative p -cycle; that is, $\partial c \in C_{p-1}(K_0)$. In K_0 , $\partial^2 c = 0$ implies that ∂c is also a cycle and therefore represents a class in $H_{p-1}(K_0)$. This defines the connecting homomorphism, mapping the relative homology class generated by $c + C_p(K_0)$ to the absolute homology class generated by ∂c . Indeed, any cycle in the same relative class with c can be written as $c + c' + c_0$, where $c' \in B_p(K)$ and $c_0 \in C_p(K_0)$. But then $\partial c' = 0$ and ∂c_0 is a boundary in K_0 . Hence, $\partial(c + c' + c_0) = \partial c + \partial c_0$, which is homologous to ∂c as a cycle in K_0 .

As an example, consider the pair $(\mathbb{B}^3, \mathbb{S}^1)$, the 3-ball modulo its equator, triangulated by (K, K_0) . We can use the exact homology sequence to figure out the relative homology of this pair. Using reduced homology, all groups of \mathbb{B}^3 are zero. Similarly, the only non-zero reduced homology group of \mathbb{S}^1 is the first one, which has rank one. Except for dimension $p = 2$, we therefore have

$$\dots \rightarrow 0 \rightarrow H_p(\mathbb{B}^3, \mathbb{S}^1) \rightarrow 0 \rightarrow \dots ,$$

implying that $H_p(\mathbb{B}^3, \mathbb{S}^1)$ itself is zero. For $p = 2$ we have the only non-trivial portion of the long exact sequence,

$$\dots \rightarrow 0 \rightarrow H_2(\mathbb{B}^3, \mathbb{S}^1) \rightarrow \tilde{H}_1(\mathbb{S}^1) \rightarrow 0 \rightarrow \dots .$$

The map between the middle two groups is thus injective as well as surjective, which implies it is an isomorphism so that $H_2(\mathbb{B}^3, \mathbb{S}^1)$ has rank one, the same as $H_1(\mathbb{S}^1)$. Indeed, we have a single non-trivial relative homology class of dimension 2, namely the one generated by the disk spanned by the equator circle. The connecting homomorphism maps this class to the absolute homology class of dimension 1 generated by the circle itself.

Bibliographic notes. Relative homology groups were introduced in the 1920s by Solomon Lefschetz for application to his fixed point theorem. They seem barely more than an afterthought to absolute homology groups. Nevertheless, they have many applications, including the study of the local homology of a space, see e.g. [82, 116], and the computation of absolute homology groups via exact sequences.

IV.4 Exact Sequences

As we have seen above, long exact sequences are handy for deriving homology groups from other homology groups. In this section, we introduce a general method for constructing such sequences and use it to get a divide-and-conquer formulation of homology, known as the Mayer-Vietoris sequence.

Chain complexes and chain maps. Freeing ourselves from the simplicial complex background, we consider a sequence of vector spaces with homomorphisms between them, $\mathcal{U} = (U_p, u_p)$ with $u_p : U_p \rightarrow U_{p-1}$. If $u_p u_{p+1} = 0$ for every p , then we call \mathcal{U} a *chain complex* and the u_p its *boundary maps*. The vanishing of the pairwise compositions of maps is all we need to define cycle groups, $Z_p(\mathcal{U}) = \ker u_p$, boundary groups, $B_p(\mathcal{U}) = \text{im } u_{p+1}$, and homology groups, $H_p(\mathcal{U}) = \ker u_p / \text{im } u_{p+1}$, in the usual way. Of course, our favorite example is the chain complex of a simplicial complex, $\mathcal{C}(K) = (\mathbb{C}_p(K), \partial_d)$.

Letting $\mathcal{V} = (V_p, v_p)$ be another chain complex, a *chain map* is a sequence of homomorphisms $\phi_p : U_p \rightarrow V_p$, one for each dimension p , that commute with the boundary maps. Specifically, $v_p \phi_p = \phi_{p-1} u_p$, for every p , but we will often drop the indices and just write $v\phi = \phi u$ to express this property. Commutativity guarantees that cycles go to cycles, $\phi_p(Z_p(\mathcal{U})) \subseteq Z_p(\mathcal{V})$, and boundaries go to boundaries, $\phi_p(B_p(\mathcal{U})) \subseteq B_p(\mathcal{V})$. Just as in the case of the induced map defined in the previous section, this implies that the chain map induces a map on homology, $(\phi_p)_* : H_p(\mathcal{U}) \rightarrow H_p(\mathcal{V})$, for every dimension p .

Letting $\mathcal{W} = (W_p, w_p)$ be a third chain complex and the $\psi_p : V_p \rightarrow W_p$ a second chain map, we call $\mathcal{U} \rightarrow \mathcal{V} \rightarrow \mathcal{W}$ *exact* at \mathcal{V} if $\ker \psi_p = \text{im } \phi_p$ for every p . A *short exact sequence* of chain complexes is a sequence of length five,

$$0 \rightarrow \mathcal{U} \xrightarrow{\phi} \mathcal{V} \xrightarrow{\psi} \mathcal{W} \rightarrow 0,$$

that begins and ends with the trivial chain complex and is exact at \mathcal{U} , \mathcal{V} , and \mathcal{W} . This means that we have a short exact sequence of vector spaces, $0 \rightarrow U_p \rightarrow V_p \rightarrow$

$W_p \rightarrow 0$, for each dimension p . Recall that this implies that each ϕ_p is injective, each ψ_p is surjective, and each V_p is isomorphic to the direct sum of U_p and W_p , although there is no natural choice for this isomorphism.

The snake or zig-zag. We are now ready to explain the general method for constructing long exact sequences of homology groups from short exact sequences of chain complexes.

SNAKE LEMMA. Let $0 \rightarrow \mathcal{U} \xrightarrow{\phi} \mathcal{V} \xrightarrow{\psi} \mathcal{W} \rightarrow 0$ be a short exact sequence of chain complexes. Then there is a well-defined map $D : H_p(\mathcal{W}) \rightarrow H_{p-1}(\mathcal{U})$, called the *connecting homomorphism*, such that

$$\dots \rightarrow H_p(\mathcal{U}) \rightarrow H_p(\mathcal{V}) \rightarrow H_p(\mathcal{W}) \xrightarrow{D} H_{p-1}(\mathcal{U}) \rightarrow \dots$$

is a long exact sequence.

Other than the connecting homomorphism, the maps in the long exact sequence are induced by the chain maps. Before looking at the algebraic details of the construction, let us see how the Snake Lemma gives rise to the exact homology sequence of a pair described in the previous section. We have a simplicial complex, K , and a subcomplex, $K_0 \subseteq K$. The inclusions of K_0 in K and K in K induce a short exact sequence of chain complexes,

$$0 \rightarrow \mathcal{C}(K_0) \rightarrow \mathcal{C}(K) \rightarrow \mathcal{C}(K, K_0) \rightarrow 0,$$

where $\mathcal{C}(K) = (\mathcal{C}_p(K), \partial_p)$ with similar notation for K_0 and (K, K_0) . Indeed, $\mathcal{C}(K_0) \rightarrow \mathcal{C}(K)$ is injective and $\mathcal{C}(K) \rightarrow \mathcal{C}(K, K_0)$ is surjective. Finally, we have exactness in the middle because a chain of K is carried by K_0 iff it is zero in (K, K_0) . The implied long exact sequence is the exact homology sequence of (K, K_0) :

$$\dots \rightarrow H_p(K_0) \rightarrow H_p(K) \rightarrow H_p(K, K_0) \xrightarrow{D} H_{p-1}(K_0) \rightarrow \dots .$$

As always, the crucial piece of the sequence is the connecting homomorphism. We now give a detailed description of its construction, in the general setting of the Snake Lemma. We omit the proof that the long exact sequence is in fact exact, leaving that as an exercise for the interested reader.

Connecting homomorphism. We construct D in four steps using the portion of the short exact sequence of chain complexes shown below. The vertical arrows are boundary maps, and the horizontal arrows are chain maps. To simplify the discussion, we will frequently suppress the subscripts that indicate the dimensions on the maps ϕ , ψ , u , v , w , or even the names of the maps themselves, as they can be determined from the domain and the clutter they introduce is more confusing

than it is helpful:

$$\begin{array}{ccccccc}
 & & V_{p+1} & \rightarrow & W_{p+1} & \rightarrow & 0 \\
 & & \downarrow & & \square_3 & \downarrow & \\
 0 & \rightarrow & U_p & \rightarrow & V_p & \rightarrow & W_p & \rightarrow & 0 \\
 & & \downarrow & & \square_2 & \downarrow & \square_0 & \downarrow & \\
 0 & \rightarrow & U_{p-1} & \rightarrow & V_{p-1} & \rightarrow & W_{p-1} & \rightarrow & 0 \\
 & & \downarrow & & \square_1 & \downarrow & & & \\
 0 & \rightarrow & U_{p-2} & \rightarrow & V_{p-2} & & & &
 \end{array}$$

Notice the labeled commutative squares in the diagram. We will refer to them when we want to emphasize that the maps around their boundaries commute. For example, the fact that \square_0 is a commutative square means that $w\psi = \psi v$ as a map from V_p to W_{p-1} . With this in mind, here are the steps in establishing the connecting homomorphism.

STEP 1: define γ . We begin with a cycle $\alpha \in W_p$ representing a class in $H_p(W)$. Since ψ is surjective, there exists a chain $\beta \in V_p$ such that $\psi(\beta) = \alpha$. Since α has zero boundary and \square_0 is commutative, the boundary of β lies in the kernel of the second chain map, $v(\beta) \in \ker \psi$. Exactness at V_{p-1} then implies that there exists a chain $\gamma \in U_{p-1}$ whose image under the first chain map is the boundary of β , $\phi(\gamma) = v(\beta)$. We summarize the situation by extracting the relevant piece of the above diagram, and a little more:

$$\begin{array}{ccccc}
 & \beta & \xrightarrow{\psi} & \alpha & \\
 & \downarrow & & \square_0 & \downarrow \\
 \gamma & \xrightarrow{\phi} & v(\beta) & \xrightarrow{\psi} & 0 \\
 & \downarrow & & \square_1 & \downarrow \\
 0 & \xrightarrow{\phi} & 0. & &
 \end{array}$$

STEP 2: show γ is a cycle. By commutativity of \square_1 and the fact that $vv = 0$, we have $\phi u(\gamma) = 0$. But ϕ is injective, so this implies that $u(\gamma) = 0$, meaning γ is a cycle; see the diagram above. Hence, γ represents a class in $H_{p-1}(U)$, and this class is the image of the class represented by α under the connecting homomorphism. The map goes left, from α to β , then down to $v(\beta)$, and then left, again, to γ . We may draw this as a snake or a zig-zag cutting through the diagram, hence the name. Notice, however, that we have made choices for α and β and we need to show that our answer does not depend on them.

STEP 3: show independence from the choice of β . Suppose first that we make another choice for β , call it β_0 , and let γ_0 be the unique element of U_{p-1} such that $\phi(\gamma_0) = v(\beta_0)$. We again summarize the situation by extracting a piece of the diagram:

$$\begin{array}{ccccc}
 \mu & \xrightarrow{\phi} & \beta, \beta_0 & \xrightarrow{\psi} & \alpha \\
 \downarrow & & \square_2 & \downarrow & \square_0 \downarrow \\
 \gamma, \gamma_0 & \xrightarrow{\phi} & v(\beta), v(\beta_0) & \xrightarrow{\psi} & 0.
 \end{array}$$

We have $\psi(\beta) = \psi(\beta_0) = \alpha$ and therefore $\beta + \beta_0 \in \ker \psi = \text{im } \phi$, so there exists a chain $\mu \in U_p$ with $\phi(\mu) = \beta + \beta_0$. Since \square_2 commutes, $\phi u(\mu) = \phi(\gamma) + \phi(\gamma_0)$. But ϕ is injective, so $u(\mu) = \gamma + \gamma_0$. In words, γ and γ_0 differ by a boundary, namely $u(\mu)$, and therefore represent the same homology class.

STEP 4: show independence from choice of α . Finally, we consider what happens with a different choice of α , say α_0 . Let β_0 and γ_0 be defined from α_0 , the same way β and γ are defined from α . Since α and α_0 are two choices of representative for the same homology class in $H_p(W)$, there exists a chain ν in W_{p+1} such that $w(\nu) = \alpha + \alpha_0$. Since ψ is surjective, there exists a chain $\varrho \in V_{p+1}$ with $\psi(\varrho) = \nu$. The situation is again summarized in a portion of the diagram:

$$\begin{array}{ccccc}
& \varrho & \xrightarrow{\psi} & \nu & \\
& \downarrow & & \square_3 & \downarrow \\
\mu' & \xrightarrow{\phi} & v(\varrho), \beta, \beta_0 & \xrightarrow{\psi} & \alpha, \alpha_0 \\
\downarrow & \square_2 & \downarrow & \square_0 & \downarrow \\
\gamma, \gamma_0 & \xrightarrow{\phi} & 0, v(\beta), v(\beta_0) & \xrightarrow{\psi} & 0.
\end{array}$$

By commutativity of \square_3 , $v(\varrho)$ and $\beta + \beta_0$ both map to $\alpha + \alpha_0$. This implies that their sum lies in $\ker \psi = \text{im } \phi$ and there is a chain μ' in U_p with $\phi(\mu') = v(\varrho) + \beta + \beta_0$. Using commutativity of \square_2 and $vv = 0$, we see that $\phi u(\mu') = v(\beta + \beta_0)$. But injectivity of ϕ implies that the preimage of $v(\beta + \beta_0)$ is $\gamma + \gamma_0$ and hence $u(\mu') = \gamma + \gamma_0$. We see that γ and γ_0 differ by a boundary and thus represent the same homology class, as required. This finishes the construction of the connecting homomorphism, D .

Mayer-Vietoris sequence. We use the Snake Lemma to derive the divide-and-conquer formulation of homology known as the Mayer-Vietoris sequence. Given two spaces, it relates their homology to the homology of their union and their intersection.

MAYER-VIETORIS SEQUENCE THEOREM. Let K be a simplicial complex and K', K'' subcomplexes such that $K = K' \cup K''$. Let $A = K' \cap K''$. Then there exists a long exact sequence

$$\dots \rightarrow H_p(A) \rightarrow H_p(K') \oplus H_p(K'') \rightarrow H_p(K) \rightarrow H_{p-1}(A) \rightarrow \dots$$

and similarly for the reduced homology groups.

PROOF. On the level of chains, $C_p(A)$ is a subgroup of both $C_p(K')$ and $C_p(K'')$. Forming the direct sums, $C_p(K') \oplus C_p(K'')$, for all dimensions p , we get a chain complex $C(K') \oplus C(K'')$ with boundary map defined componentwise. We have two copies of $C_p(A)$ and can kill one off with the image of $C_p(A)$ via the diagonal, and the quotient is easily identified with $C_p(K)$. Stated more formally, let $i' : A \rightarrow K'$ and $i'' : A \rightarrow K''$ be the inclusions of A , and let $j' : K' \rightarrow K$ and $j'' : K'' \rightarrow K$

be the inclusions into K . Set $i(a) = (i'(a), i''(a))$ and $j(x, y) = j'(x) + j''(y)$. Then it is not difficult to see that we have a short exact sequence of chain complexes, namely

$$0 \rightarrow \mathcal{C}(A) \xrightarrow{i} \mathcal{C}(K') \oplus \mathcal{C}(K'') \xrightarrow{j} \mathcal{C}(K) \rightarrow 0.$$

The long exact sequence implied by the Snake Lemma is the Mayer-Vietoris sequence. The above is easily adapted to the reduced sequence as well. \square

Exactness of the Mayer-Vietoris sequence at $H_p(K)$ tells us that this group is isomorphic to the direct sum of the image of $j_* : H_p(K') \oplus H_p(K'') \rightarrow H_p(K)$ with the kernel of $i_* : H_{p-1}(A) \rightarrow H_{p-1}(K') \oplus H_{p-1}(K'')$. This distinguishes two types of homology classes in K . A class in $\text{im } j_*$ lives in K' , in K'' , or in both. A class in $\ker i_*$ corresponds to a $(p-1)$ -dimensional cycle $\gamma \in A$ that bounds both in K' and K'' . If we write $\gamma = \partial\alpha' = \partial\alpha''$, with α' a p -chain in K' and α'' a p -chain in K'' , then $\alpha = \alpha' + \alpha''$ is a cycle in K that represents this second type of class; see Figure IV.10. It is useful to check through the four steps constructing

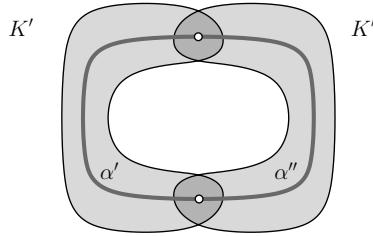


Figure IV.10: The 1-cycle α decomposes into 1-chains α' in K' and α'' in K'' . The common boundary of the two 1-chains is a pair of points, a reduced 0-cycle in A .

the connecting homomorphism, D . They take a class in $H_p(K)$ and define one in $H_{p-1}(A)$ as follows. Represent the class by α , a p -cycle of K . As before, there exists β , a p -chain in $C_p(K') \oplus C_p(K'')$, such that $j(\beta) = \alpha$. In fact, there are several, and we get them by writing $\alpha = \alpha' + \alpha''$, with α' in K' , α'' in K'' , and setting $\beta = (\alpha', \alpha'')$. Different ways of decomposing α give different β , but note that any two differ by something in A . Now take $\partial\beta = (\partial\alpha', \partial\alpha'')$. The fact that α is a cycle tells us that $\partial\alpha' = \partial\alpha''$ and that it lies in A . Thus, the cycle γ in the construction of D is $\partial\alpha'$.

The sphere, \mathbb{S}^d . To illustrate the utility of the Mayer-Vietoris sequence, we use it to compute the homology of the d -dimensional sphere, \mathbb{S}^d . Specifically, we show that

$$\tilde{\beta}_p(\mathbb{S}^d) = \begin{cases} 1 & \text{if } p = d; \\ 0 & \text{if } p \neq d. \end{cases}$$

Writing the sphere as the union of its upper and lower hemisphere, $\mathbb{S}^d = U \cup L$, we get the equator as the intersection, $A = U \cap L$. Each hemisphere is a ball, and the

equator is a sphere of dimension $d - 1$. This allows us to compute the homology of \mathbb{S}^d inductively, using the reduced Mayer-Vietoris sequence:

$$\dots \rightarrow \tilde{\mathcal{H}}_p(A) \rightarrow \tilde{\mathcal{H}}_p(U) \oplus \tilde{\mathcal{H}}_p(L) \rightarrow \tilde{\mathcal{H}}_p(\mathbb{S}^d) \rightarrow \tilde{\mathcal{H}}_{p-1}(E) \rightarrow \dots .$$

For $d = 0$, the sphere is two points, so its reduced homology has rank one in dimension 0 and rank zero otherwise. This established the induction basis. For general d , the sequence decomposes into pieces of the form

$$0 \oplus 0 \rightarrow \tilde{\mathcal{H}}_p(\mathbb{S}^d) \rightarrow \tilde{\mathcal{H}}_{p-1}(\mathbb{S}^{d-1}) \rightarrow 0 \oplus 0,$$

where $0 \oplus 0$ is of course the zero element in the direct sum of the homology groups of the two hemispheres. This implies that the rank of the p -th reduced homology group of \mathbb{S}^d is the same as the rank of the $(p-1)$ -st reduced homology group of \mathbb{S}^{d-1} , namely one for $p = d$ and zero otherwise, as claimed. Note that the generator of $\tilde{\mathcal{H}}_d(\mathbb{S}^d)$ is the second type of class, consisting of two chains, one from each hemisphere, whose boundary is the generating cycle of $\tilde{\mathcal{H}}_{d-1}(\mathbb{S}^{d-1})$. In particular, it is represented by the sum of all its d -dimensional simplices.

The real projective space, \mathbb{P}^d . As another example, we consider the real projective d -dimensional space which is the quotient space of \mathbb{S}^d by the antipodal map, $f(x) = -x$. In other words, \mathbb{P}^d is obtained by gluing \mathbb{S}^d to itself by identifying antipodal points in pairs. We show that its reduced Betti numbers are

$$\tilde{\beta}_p(\mathbb{P}^d) = \begin{cases} 1 & \text{for } 1 \leq p \leq d; \\ 0 & \text{otherwise.} \end{cases}$$

For dimensions $d = 0, 1$ we have familiar spaces, namely the point, \mathbb{P}^0 , and the circle, \mathbb{P}^1 . We already know their homology and their reduced Betti numbers agree with the claimed formula. This establishes the induction basis. For general d , we decompose \mathbb{S}^d into three subspaces by limiting the d -th coordinate to $x_d \leq -1/2$, $-1/2 \leq x_d \leq 1/2$, and $1/2 \leq x_d$. The first and the last are identified by f and give a single subspace $B \subseteq \mathbb{P}^d$, which is a ball. The middle subspace becomes a space M that is homotopy equivalent to the quotient space of the equator, where $x_d = 0$, which is in turn homeomorphic to \mathbb{P}^{d-1} . The middle subspace intersects the union of the other two in two spheres of dimension $d - 1$. Taking the quotient identifies the two spheres, implying that B and M intersect in a single sphere of dimension $d - 1$. Since the reduced homology of B vanishes in all dimensions p , the Mayer-Vietoris sequence decomposes into pieces of the form

$$0 \rightarrow 0 \oplus \tilde{\mathcal{H}}_p(\mathbb{P}^{d-1}) \rightarrow \tilde{\mathcal{H}}_p(\mathbb{P}^d) \rightarrow 0,$$

for $p < d - 1$. By induction, this establishes $\tilde{\beta}_0(\mathbb{P}^d) = 0$ and $\tilde{\beta}_p(\mathbb{P}^d) = 1$ for $1 \leq p \leq d - 1$. We still need to show that the d -th reduced Betti number is equal to one. The piece of the Mayer-Vietoris sequence we use for this is

$$0 \oplus 0 \rightarrow \tilde{\mathcal{H}}_d(\mathbb{P}^d) \xrightarrow{D} \tilde{\mathcal{H}}_{d-1}(\mathbb{S}^{d-1}) \xrightarrow{g_*} 0 \oplus \tilde{\mathcal{H}}_{d-1}(\mathbb{P}^{d-1}) \rightarrow \tilde{\mathcal{H}}_{d-1}(\mathbb{P}^d) \rightarrow 0.$$

We claim that the map g_* is 0. This implies that D is surjective, and since it is also injective, $\tilde{\beta}_d(\mathbb{P}^d) = 1$, as required. To see that g_* is zero, we use the inductive assumption, namely that $\tilde{\beta}_{d-1}(\mathbb{P}^{d-1}) = 1$. The corresponding homology group has a unique generator, namely the sum of all $(d-1)$ -simplices triangulating \mathbb{P}^{d-1} . The map g takes each simplex in the triangulation of \mathbb{S}^{d-1} to its quotient, which means each simplex in \mathbb{P}^{d-1} is counted twice. The top-dimensional simplices cancel in pairs, which completes the calculation.

Bibliographic notes. The introduction of exact sequences is often attributed to Eilenberg and sometimes to Lyndon, but see also [90]. The Snake Lemma is a major achievement of algebraic topology, and the construction of the connecting homomorphism is its critical piece. A complete proof can be found in many algebraic topology texts, including [116]. The Mayer-Vietoris sequences are older than the Snake Lemma and go back to the work by Mayer [108] and Vietoris [147].

Exercises

The credit assignment reflects a subjective assessment of difficulty. A typical question can be answered using knowledge of the material combined with some thought and analysis.

1. **Sperner Lemma** (three credits). Let K be a triangulated triangular region as in Figure IV.11. We 3-color the vertices such that

- the three corners receive three different colors;
- the vertices on each side of the region are 2-colored.

Prove that there is a triangle in K whose vertices receive three different colors.

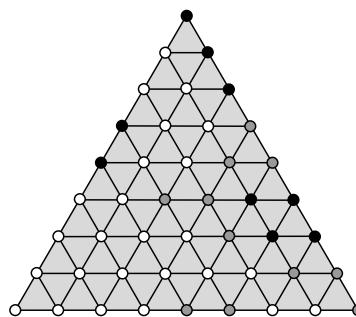


Figure IV.11: Each vertex receives one of three colors: white, shaded, or black.

2. **Isomorphic homology** (one credit). Construct two topological spaces that have isomorphic homology groups but are not homotopy equivalent.

3. **Fixed point** (two credits). Let $f : \mathbb{B}^d \rightarrow \mathbb{B}^d$ be a continuous map with the property that there is a $\delta < 1$ such that $\|f(x) - f(y)\| \leq \delta \|x - y\|$ for all points $x, y \in \mathbb{B}^d$. In words, the distance between any two points diminishes by at least a constant factor $\delta < 1$ each time we apply f . Prove that such a map f has a unique fixed point $x = f(x)$. [On orientation maps, this point is usually marked as “you are here”].
4. **Klein bottle** (one credit). Show that the Betti numbers of the 2-dimensional Klein bottle are $\beta_0 = 1$, $\beta_1 = 2$, $\beta_2 = 1$. Which other 2-manifold has the same Betti numbers?
5. **Dunce cap** (three credits). The *dunce cap* is constructed from a piece of cloth in the shape of an equilateral triangle as follows. Orienting two edges away from a common origin, we glue them to each other as prescribed by their orientation. This gives a piece of a cone with a rim (the third edge) and a seam (the glued first two edges). Now we orient the rim and glue it along the seam, again such that orientations match. The result reminds us of the shell of a snail, perhaps.
 - (i) Give a triangulation of the dunce cap.
 - (ii) Show that the reduced Betti numbers of the dunce cap vanish in all dimensions.
 - (iii) Show that the dunce cap is contractible but any triangulation of it is not collapsible.
6. **3-torus** (three credits). Consider the *3-dimensional torus* obtained from the unit cube by gluing opposite faces in pairs, without twisting. That is, each point $(x, y, 0)$ is identified with $(x, y, 1)$, $(x, 0, z)$ with $(x, 1, z)$, and $(0, y, z)$ with $(1, y, z)$. Show that the Betti numbers of this space are $\beta_0 = \beta_3 = 1$ and $\beta_1 = \beta_2 = 3$.
7. **The Steenrod Five Lemma** (two credits). Suppose we have a commutative diagram of vector spaces and homomorphisms,

$$\begin{array}{ccccccc} U_1 & \rightarrow & U_2 & \rightarrow & U_3 & \rightarrow & U_4 & \rightarrow & U_5 \\ \downarrow & & \downarrow & & \downarrow & & \downarrow & & \downarrow \\ V_1 & \rightarrow & V_2 & \rightarrow & V_3 & \rightarrow & V_4 & \rightarrow & V_5, \end{array}$$

where the horizontal sequences are exact at the middle three vector spaces and the first two and last two vertical arrows are isomorphisms. Prove that the middle vertical arrow is then also an isomorphism.

8. **Exact sequence of a triple** (one credit). Let C be a simplicial complex with subcomplexes $A \subseteq B \subseteq C$. Prove the existence of the following *exact homology sequence of the triple*:

$$\dots \rightarrow H_p(B, C) \rightarrow H_p(A, C) \rightarrow H_p(A, B) \rightarrow H_{p-1}(B, C) \rightarrow \dots .$$

Chapter V

Duality

Instead of computing homology from a triangulation, we can also work with different decompositions and get isomorphic groups. The alpha complex and the dual Voronoi decomposition of a union of balls come to mind. Generalizing this geometric idea beyond Euclidean space, and in particular beyond manifolds, runs into difficulty. This is the motivation for taking the issue to the algebraic level, where it leads to the concept of cohomology groups. For modulo 2 arithmetic, these are isomorphic to the corresponding homology groups, but the isomorphisms are not natural. For nice topological spaces, such as manifolds and manifolds with boundary, there are relations between the homology and the cohomology groups that go beyond the general relations. In this chapter, we will see three of them: Poincaré duality, Lefschetz duality, and Alexander duality. The last of the three has algorithmic ramifications for subsets of 3-dimensional Euclidean space.

V.1 Cohomology

In this section, we introduce cohomology groups. They are similar to homology groups but less geometric and motivated primarily by algebraic considerations. They belong to the standard tool set of an algebraic topologist and appear in modern statements of the duality results discussed in the subsequent three sections.

Groups of maps. Let $G = \mathbb{Z}_2$, the group of two elements, 0 and 1, together with addition modulo 2. All abelian groups we have encountered so far are vector spaces isomorphic to G^n for some integer n . Let U be such a vector space and let $\varphi : U \rightarrow G$ be a homomorphism. To define φ , it suffices to specify its values on the generators of U . If φ_0 is a second such homomorphism, their sum is defined by $(\varphi + \varphi_0)(u) = \varphi(u) + \varphi_0(u)$. This is again a homomorphism because

$$\begin{aligned} (\varphi + \varphi_0)(u + v) &= \varphi(u + v) + \varphi_0(u + v) \\ &= (\varphi + \varphi_0)(u) + (\varphi + \varphi_0)(v), \end{aligned}$$

It is easy to see that addition of homomorphisms is associative. We also have a neutral element, the zero homomorphism that sends every $u \in U$ to $0 \in G$, and an inverse, which for modulo 2 arithmetic is the identity, $-\varphi = \varphi$. We thus have a *group of homomorphisms* from U to G , denoted as $\text{Hom}(U, G)$. Think for example of U as the group of p -chains of a simplicial complex and $\text{Hom}(U, G)$ as the group of labelings of the p -simplices by 0 and 1. The vector spaces U and $\text{Hom}(U, G)$ are isomorphic, although the isomorphism requires us to pick a basis of U . Specifically, if U has the basis e_1, e_2, \dots, e_n , then $\text{Hom}(U, G)$ has the basis f_1, f_2, \dots, f_n , where $f_i(e_i) = 1$ and $f_i(e_j) = 0$ whenever $i \neq j$. The corresponding isomorphism is defined by mapping e_i to f_i for all i . If we choose a different basis, the isomorphism changes.

We give a specific example to emphasize this point. Take $U = G^2$, with first basis $e_1 = (1, 0)$ and $e_2 = (0, 1)$. An element $w = (a, b)$ in U is written as $ae_1 + be_2$. The isomorphism from U to $\text{Hom}(U, G)$ thus takes w to the map $f_w = af_1 + bf_2$. Suppose instead that we consider the basis $e'_1 = (1, 1) = e_1 + e_2$ and $e'_2 = (0, 1) = e_2$. Then $w = (a, b)$ is written as $ae'_1 + (b - a)e'_2$. The new isomorphism takes w to the map $f'_w = af'_1 + (b - a)f'_2$. To see that f_w and f'_w are generally different, we evaluate both at $v = (x, y)$ in U , giving

$$\begin{aligned} f_w(v) &= af_1(v) + bf_2(v) &= ax + by, \\ f'_w(v) &= af'_1(v) + (b - a)f'_2(v) &= ax + (b - a)(y - x). \end{aligned}$$

The two isomorphisms are thus indeed different. This is the reason for why cohomology is worth defining at all, because if there were a natural isomorphism between U and $\text{Hom}(U, G)$, the theories of homology and cohomology would be the same.

Given another vector space V and a homomorphism $f : U \rightarrow V$, there is an induced *dual homomorphism*, $f^* : \text{Hom}(V, G) \rightarrow \text{Hom}(U, G)$, that maps $\psi : V \rightarrow G$ to the composite $f^*(\psi) = \psi \circ f : U \rightarrow G$. The map f^* is indeed a homomorphism since

$$\begin{aligned} f^*(\psi + \psi_0)(u) &= (\psi + \psi_0) \circ f(u) \\ &= \psi \circ f(u) + \psi_0 \circ f(u) \\ &= f^*(\psi)(u) + f^*(\psi_0)(u) \end{aligned}$$

for every $u \in U$. The group of homomorphisms and the dual homomorphism can be defined for more general abelian groups U , V , and G , but this will not be necessary for our purposes.

Simplicial cohomology. Let K be a simplicial complex. We construct cohomology groups by turning chain groups into groups of homomorphisms and boundary maps into their dual homomorphisms. To begin, we define a *p -cochain* as a homomorphism $\varphi : C_p \rightarrow G$, where $G = \mathbb{Z}_2$ as before. Given a p -chain, $c \in C_p$, the cochain evaluates c by mapping it to 0 or 1. It is common to write this evaluation like a scalar product, $\varphi(c) = \langle \varphi, c \rangle$. Letting ℓ be the number of p -simplices σ in c with $\varphi(\sigma) = 1$, we have $\langle \varphi, c \rangle = 1$ iff ℓ is odd. Considering chains and cochains as sets, the evaluation thus distinguishes odd from even numbers of intersections.

The p -dimensional cochains form the *group of p -cochains*, $C^p = \text{Hom}(C_p, G)$. Recall that the boundary map is a homomorphism $\partial_p : C_p \rightarrow C_{p-1}$. It thus defines a dual homomorphism, the *coboundary map*

$$\delta^{p-1} : \text{Hom}(C_{p-1}, G) \rightarrow \text{Hom}(C_p, G),$$

or simply $\delta : C^{p-1} \rightarrow C^p$. It is worth looking at this construction in more detail. Let φ be a $(p-1)$ -cochain and ∂c a $(p-1)$ -chain. By definition of dual homomorphism, φ applied to ∂c is the same as $\delta\varphi$ applied to c ; that is, $\langle \varphi, \partial c \rangle = \langle \delta\varphi, c \rangle$. Suppose for example that φ evaluates a single $(p-1)$ -simplex to one and all others to zero. Then $\delta\varphi$ evaluates all p -dimensional cofaces of this simplex to one and all others to zero. This gives a concrete interpretation of the coboundary map, which will allow us to construct more elaborate examples shortly. Since the coboundary map runs in a direction opposite to the boundary map, it raises the dimension. Its kernel is the *group of cocycles*, and its image is the *group of coboundaries*:

$$\begin{aligned} Z^p &= \ker \delta^p : C^p \rightarrow C^{p+1}, \\ B^p &= \text{im } \delta^{p-1} : C^{p-1} \rightarrow C^p. \end{aligned}$$

Recall the Fundamental Lemma of Homology according to which $\partial \circ \partial : C_{p+1} \rightarrow C_{p-1}$ is the zero homomorphism. We therefore have $\langle \delta\delta\varphi, c \rangle = \langle \delta\varphi, \partial c \rangle = \langle \varphi, \partial\partial c \rangle = 0$. In other words, $\delta \circ \delta : C^{p-1} \rightarrow C^{p+1}$ is also the zero homomorphism. Hence, the coboundary groups are subgroups of the cocycle groups, and we have the familiar picture, except that the maps now go from right to left, as in Figure V.1.

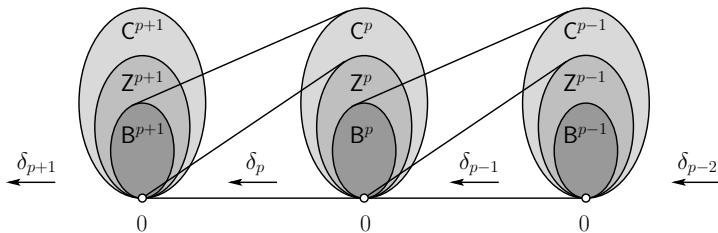


Figure V.1: The cochain complex consisting of a linear sequence of cochain, cocycle, and coboundary groups connected by coboundary homomorphisms.

DEFINITION. The p -th *cohomology group* is the quotient of the p -th cocycle group modulo the p -th coboundary group, $H^p = Z^p / B^p$, for all p .

Reduced cohomology. Similar to homology, it is often useful to modify the definition slightly and to define the *reduced cohomology groups*, denoted as \tilde{H}^p . Recall that for homology, this is done by introducing the augmentation map $\epsilon : C_0 \rightarrow \mathbb{Z}_2$ defined by $\epsilon(u) = 1$ for each vertex u . The (-1) -st cochain group, $C^{-1} = \text{Hom}(\mathbb{Z}_2, G)$, has two elements, the map ϕ_0 mapping 1 to 0 and the map ϕ_1 mapping 1 to 1. The dual homomorphism of the augmentation map, $\epsilon^* : \text{Hom}(\mathbb{Z}_2, G) \rightarrow C^0$,

maps ϕ_0 to ψ_0 , which evaluates every vertex to zero, and ϕ_1 to ψ_1 , which evaluates every vertex to one. With this, we have

$$\dots \xleftarrow{\delta^1} C^1 \xleftarrow{\delta^0} C^0 \xleftarrow{\epsilon^*} \text{Hom}(\mathbb{Z}_2, G) \xleftarrow{0} 0 \xleftarrow{0} \dots.$$

Before the modification, the only 0-coboundary was the trivial 0-cochain, ψ_0 . Now we have two 0-coboundaries, ψ_0 and ψ_1 . The net effect of this modification is that the rank of the zeroth cohomology group drops by one, similar to the rank of the zeroth homology group when we add the augmentation map. As an exception to this rule, the ranks of H^0 and \tilde{H}^0 are the same if C^0 is trivial, in which case $\text{rank } \tilde{H}^{-1} = 1$, similar to reduced homology.

An example. To get a better feeling for cohomology, let us consider the triangulation of the annulus in Figure V.2. The 0-cochain that evaluates every vertex to one is a 0-cocycle because every edge has exactly two vertices, which implies that the coboundary of this particular 0-cochain is the zero homomorphism. This is the only non-trivial 0-cocycle, and since for dimensional reasons there are no non-trivial 0-coboundaries, this implies that the zeroth cohomology group, H^0 , has rank one. Correspondingly, the zeroth reduced cohomology group is zero.

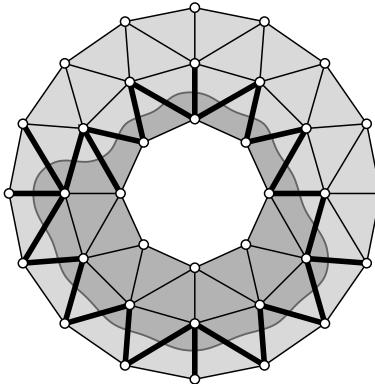


Figure V.2: The 1-cocycle is drawn by highlighting the edges it evaluates to one. They all cross the “dual” closed curve. The 1-cocycle is a 1-coboundary because it is the coboundary of the 0-cochain that evaluates a vertex to one iff it lies in the shaded region inside the closed curve.

One dimension up, we consider a 1-cochain $\varphi : C_1 \rightarrow G$. Its coboundary is the 2-chain $\delta\varphi : C_2 \rightarrow G$ that evaluates a triangle to one iff it is the coface of an odd number of edges evaluated to one by φ . Hence, φ is a 1-cocycle iff every triangle is incident to an even number of edges evaluating to one. A 1-cocycle thus looks like a picket fence; see Figure V.2. In this example, we can draw a closed curve such that an edge evaluates to one iff it crosses the curve. It follows that a 1-chain is evaluated to the parity of the number of times it crosses that curve. If the 1-chain is a 1-cycle, then this number is necessarily even and the evaluation is zero. The 1-cocycle in

Figure V.2 is also a 1-coboundary. To get a 1-cocycle that is not the image of a 0-cochain, we construct a picket fence that starts with an outer boundary edge of the annulus and ends with an inner boundary edge. All such picket fences are cohomologous, and any one of them can be used as representative of the cohomology class that generates the first cohomology group. The rank of H^1 is therefore one.

Another dimension up, we have $Z^2 = C^2$ simply because every 2-cochain maps to zero, the sole element of C^3 . We also have $B^2 = C^2$. To see this, note that the 2-cochain that evaluates a single triangle to one and all others to zero is a 2-coboundary. Indeed, we can draw three curves from a point in the interior of the triangle to the boundary of the annulus and get a “dual” 1-cochain as the sum of three picket fences, one for each curve, whose coboundary is the 2-cochain. Other 2-cochains are obtained as coboundaries of sums of such triplets of picket fences. It follows that the second cohomology group, H^2 , has rank 0.

Observe that the ranks of the cohomology groups are the same as the ranks of the corresponding homology groups. This is not a coincidence.

Coboundary matrix. Recall that we can get the rank of the p -th homology group from two boundary matrices transformed into normal form by row and column operations. Recall also that $\text{rank } H_p = \text{rank } Z_p - \text{rank } B_p$. As illustrated in Figure V.3, the right-hand side of this equation is the number of zero columns in the p -th matrix minus the number of non-zero rows in the $(p+1)$ -st matrix; compare with Figure IV.5. As we have seen earlier, a cochain evaluates a single p -simplex to one

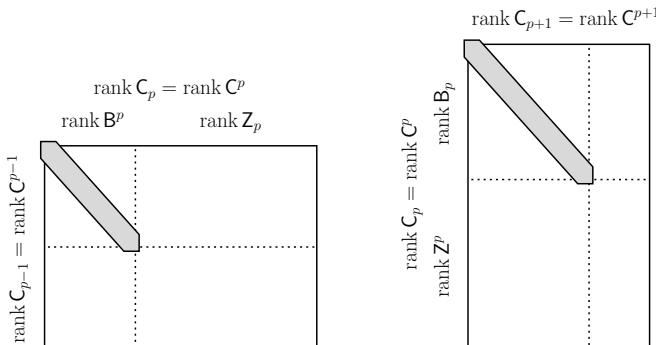


Figure V.3: The p -th and $(p+1)$ -st boundary matrices in normal form. They are also the coboundary matrices in normal form transposed.

and all others to zero iff its coboundary evaluates each $(p+1)$ -coface of this p -simplex to one and all other $(p+1)$ -simplices to zero. It follows that the coboundary matrices are the boundary matrices transposed. The normal form of the boundary matrices thus already contains the information we need to get at the ranks of the cohomology groups. Specifically, $\text{rank } H^p = \text{rank } Z^p - \text{rank } B^p$; the rank of the cocycle group is the number of zero rows in the $(p+1)$ -st boundary matrix, and the rank of the coboundary group is the number of non-zero columns in the p -th boundary matrix,

both in normal form. The number of columns of the p -th matrix is the number of rows of the $(p+1)$ -st matrix, and hence $\text{rank } \mathbf{B}^p + \text{rank } \mathbf{Z}_p = \text{rank } \mathbf{Z}^p + \text{rank } \mathbf{B}_p$; see Figure V.3. This implies

$$\begin{aligned}\text{rank } \mathbf{H}^p &= \text{rank } \mathbf{Z}^p - \text{rank } \mathbf{B}^p \\ &= \text{rank } \mathbf{Z}_p - \text{rank } \mathbf{B}_p = \text{rank } \mathbf{H}_p.\end{aligned}$$

Since homology and cohomology groups have the same rank, there is no concept of co-Betti number. For modulo 2 arithmetic, the rank determines the group; hence homology and cohomology groups are isomorphic: $\mathbf{H}_p \simeq \mathbf{H}^p$ for all p . This is the \mathbb{Z}_2 -version of a standard result in algebraic topology. For more general coefficient groups, it relates the free parts and torsion parts of the homology groups with those of the cohomology groups. A more complete statement of the result for \mathbb{Z}_2 -coefficients is the following.

UNIVERSAL COEFFICIENT THEOREM. Given a topological space, \mathbb{X} , there are maps $\mathbf{H}^p(\mathbb{X}) \rightarrow \text{Hom}(\mathbf{H}_p(\mathbb{X}), \mathbf{G}) \rightarrow \mathbf{H}_p(\mathbb{X})$ in which the first map is a natural isomorphism and the second is an isomorphism that is not natural.

We saw at the beginning of this section that the second isomorphism depends on a choice of basis and is therefore not natural. The first isomorphism does not depend on such a choice. It is natural in the sense that if \mathbb{Y} is another topological space and $f : \mathbb{X} \rightarrow \mathbb{Y}$ is a continuous map, then the diagram

$$\begin{array}{ccc}\mathbf{H}^p(\mathbb{X}) & \rightarrow & \text{Hom}(\mathbf{H}_p(\mathbb{X}), \mathbf{G}) \\ \uparrow & & \uparrow \\ \mathbf{H}^p(\mathbb{Y}) & \rightarrow & \text{Hom}(\mathbf{H}_p(\mathbb{Y}), \mathbf{G})\end{array}$$

of induced maps commutes. The fact that the isomorphism between \mathbf{H}^p and $\text{Hom}(\mathbf{H}_p, \mathbf{G})$ is natural is the reason for why there is no need to introduce a theory of co-cohomology.

Bibliographic notes. Similar to homology, cohomology is an established topic within algebraic topology today, but it took some time to become clearly established. Cohomology has a long and complicated history with a variety of precursors that go back to Poincaré, Alexander, Lefschetz, de Rham, Pontryagin, Kolmogorov, Whitney, Čech, Eilenberg, Steenrod, Spanier, and others. All these approaches were unified with the clear statement of a set of axioms that characterize homology and cohomology theories [66]. The Universal Coefficient Theorem and the duality theorems in the coming three sections were originally proven in more elementary forms before being reformulated in terms of homology and cohomology as we describe them here [116].

V.2 Poincaré Duality

For sufficiently nice topological spaces, there are relations between the homology and cohomology groups that go beyond the ones we have already seen. These

relationships go under the name of duality. The first and most important of these is Poincaré duality, which we describe in this and the next section.

Combinatorial manifolds. In the rest of this chapter, we work only with triangulations of manifolds that satisfy a condition on the topology of the links. Specifically, a *combinatorial d-manifold* is a manifold of dimension d together with a triangulation such that the link of every i -simplex triangulates the sphere of dimension $d - i - 1$. The condition implies that the closed star of every simplex has the topology of the d -dimensional ball, \mathbb{B}^d . To describe this in greater detail, we introduce the *join* of two topological spaces, \mathbb{X} and \mathbb{Y} , which we denote as $\mathbb{X} * \mathbb{Y}$. Begin with the product $\mathbb{X} \times [0, 1] \times \mathbb{Y}$. For each $x_0 \in \mathbb{X}$ identify all points $(x_0, 0, y)$ together, and for each $y_0 \in \mathbb{Y}$ identify all points $(x, 1, y_0)$ together. The quotient space of these identifications is $\mathbb{X} * \mathbb{Y}$. Also, $\mathbb{X} * \emptyset = \mathbb{X}$. Figure V.4 illustrates the construction by showing the *suspension* of \mathbb{X} , that is, the join with the 0-sphere, denoted as $\Sigma \mathbb{X} = \mathbb{X} * \mathbb{S}^0$. Geometrically, we can think of the join as a union of all line segments connecting \mathbb{X} to \mathbb{Y} , which are kept disjoint except at their endpoint in \mathbb{X} and \mathbb{Y} .

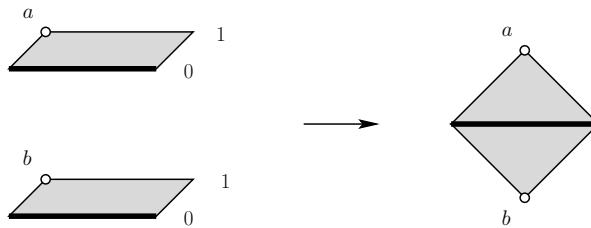


Figure V.4: Constructing the join of a line segment and a pair of points. Left: the product of the two with the unit interval. Right: the suspension obtained from the product by identification.

Returning to the definition of a combinatorial manifold, we recall that the star of a simplex, σ , consists of all simplices τ that contain σ as a face. Besides σ , each simplex in the star is the join of σ with a simplex in the link of σ . If σ is an i -simplex, then $\text{Lk } \sigma$ is a $(d - i - 1)$ -sphere. Taking the join, we get a d -ball, as mentioned earlier.

Exotic manifolds. Not every triangulation of a manifold satisfies the conditions on the links given above. We describe the construction of a triangulation of the 5-sphere that has a vertex whose link is not a 4-sphere. We begin with a triangulation, P , of the Poincaré homology 3-sphere. This space is homologically the same as but topologically different from the 3-sphere, \mathbb{S}^3 . There are many ways to describe it. A particularly convenient way uses three complex numbers to write a point in \mathbb{R}^6 . Letting x_1, x_2, \dots, x_6 be the coordinates, we set $x = x_1 + ix_2$, $y = x_3 + ix_4$, $z = x_5 + ix_6$ and recall that their conjugates are $\bar{x} = x_1 - ix_2$, $\bar{y} = x_3 - ix_4$,

$\bar{z} = x_5 - ix_6$. Consider the following two equations:

$$\begin{aligned} x\bar{x} + y\bar{y} + z\bar{z} &= 1, \\ x^2 + y^3 + z^5 &= 0. \end{aligned}$$

The first equation describes the 5-sphere. The second equation is really two equations, one for the real and the other for the imaginary parts, and it defines a 4-dimensional space whose points have neighborhoods homeomorphic to \mathbb{R}^4 except at the origin, where the space is singular. The intersection of the two spaces is the Poincaré homology 3-sphere. It is triangulable, and we let P be a triangulation of this space. Next, we take two suspension steps to construct a triangulation of the 5-sphere. Writing this in terms of triangulations, we get

$$\begin{aligned} \Sigma P &= \{a, b\} \cup \{\sigma, a * \sigma, b * \sigma \mid \sigma \in P\}, \\ \Sigma^2 P &= \{u, v\} \cup \{\tau, u * \tau, v * \tau \mid \tau \in \Sigma P\}. \end{aligned}$$

The shared link of the vertices a and b in ΣP is P , which is not a triangulation of \mathbb{S}^3 . It follows that a and b do not have neighborhoods homeomorphic to \mathbb{R}^3 . Hence, the underlying space of ΣP is not even a manifold. Taking the suspension twice is the same as forming the join with a circle. Hence, $\Sigma^2 P$ triangulates the join of the Poincaré homology 3-sphere with \mathbb{S}^1 . As it turns out, this join is homeomorphic to \mathbb{S}^5 . The proof of this fact is not easy and is omitted. But now we have a triangulation of a 5-manifold, namely $\Sigma^2 P$, that violates the condition on the links. Specifically, the shared link of the vertices u and v in $\Sigma^2 P$ is ΣP , which is not even a 4-manifold.

Dual blocks. Now let \mathbb{M} be a compact, combinatorial d -manifold triangulated by K . Recall that the barycentric subdivision, SdK , is obtained by connecting the barycenters of the simplices in K ; see Section III.1 for the definition and Figure V.5 for an illustration. It is not difficult to show that if K has the link property required for a combinatorial manifold, then so does SdK . Label each vertex in SdK by the dimension of the corresponding simplex in K and note that each simplex in SdK has distinct labels on its vertices. The vertex with smallest label is therefore unique. Letting u be the barycenter of σ in K , the *dual block*, denoted by $\hat{\sigma}$, is the union of the simplices in the barycentric subdivision for which u is the vertex with minimum label; again see Figure V.5. We let B be the set of dual blocks and call it the *dual block decomposition* of \mathbb{M} . For example, in the case of a combinatorial 3-manifold, the dual blocks to a vertex, edge, triangle, and tetrahedron are, respectively, a ball, a disk, an interval, and a point. The relationship between K and B is much like that between the Delaunay triangulation and its dual Voronoi diagram. In particular, if the p -simplex σ is a face of the $(p+1)$ -simplex τ , then the dual block $\hat{\sigma}$ contains $\hat{\tau}$ in its boundary. In fact, the boundary of $\hat{\sigma}$ is the union of dual blocks $\hat{\tau}$ over all proper cofaces τ of σ . We denote this boundary by $bd\hat{\sigma}$, noting that $\hat{\sigma}$ is the join of $bd\hat{\sigma}$ with the barycenter of σ . Since we have a combinatorial manifold, $bd\hat{\sigma}$ has the topology of the $(q-1)$ -sphere, where $p+q=d$.

We construct a new chain complex from the dual block decomposition as follows. Choosing complementary dimensions $p+q=d$, a *block chain* of dimension q is a

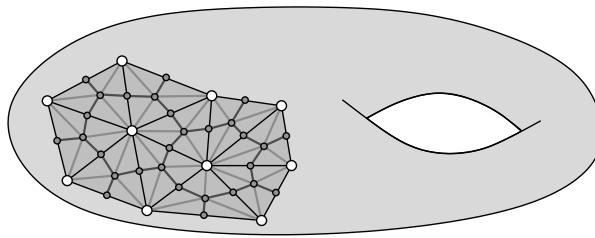


Figure V.5: A small piece of a triangulation of the torus, the barycentric subdivision, and the dual block decomposition.

formal sum $\sum a_i \hat{\sigma}_i$, where the σ_i are the p -simplices of K and the $\hat{\sigma}_i$ are the dual blocks of dimension q , with modulo 2 coefficients as usual. The collection of block chains of dimension q forms an abelian group, D_q . The boundary homomorphism connecting the q -th group to the $(q-1)$ -st group is defined by mapping $\hat{\sigma}_i$ to $\partial_q \hat{\sigma}_i = \sum \hat{\tau}_j$, where the sum is over all $(p+1)$ -dimensional cofaces τ_j of σ_i . The full boundary homomorphism, $\partial_q : D_q \rightarrow D_{q-1}$, is the linear extension to block chains. It is easy to see that $\partial_{q-1} \circ \partial_q = 0$, so that (D_q, ∂_q) is indeed a chain complex.

Blocks or simplices. We now have three ways to compute the homology of M : using the simplices in K , using the simplices in SdK , or using the dual blocks in B . We formally prove what is to be expected, namely that SdK and B give the same homology. Write $C = (C_p, \partial_p)$ for the chain complex defined by SdK and $D = (D_p, \partial_p)$ for the chain complex defined by B . Mapping each p -dimensional dual block to the sum of p -simplices it contains, we get a homomorphism $b_p : D_p \rightarrow C_p$. The maps b_p commute with the boundary maps and thus form a chain map between the two chain complexes, which we denote as $b : D \rightarrow C$.

BLOCK COMPLEX LEMMA. The chain map $b : D \rightarrow C$ induces $b_* : H_p(D) \rightarrow H_p(C)$, which is an isomorphism for each dimension p .

PROOF. Let X_p be the subcomplex of SdK consisting of all simplices that lie in blocks of dimension at most p . Clearly, $X_0 \subseteq X_1 \subseteq \dots \subseteq X_d = SdK$. The p -th relative homology group of the pair (X_p, X_{p-1}) is isomorphic to D_p . More generally,

$$H_p(X_q, X_{q-1}) \simeq \begin{cases} D_p & \text{if } p = q; \\ 0 & \text{if } p \neq q. \end{cases}$$

Indeed, each pair $(\hat{\sigma}, \text{bd } \hat{\sigma})$ has the homology of a ball relative to its boundary. Next, consider the long exact sequence of the pair (X_q, X_{q-1}) :

$$\dots \rightarrow H_{p+1}(X_q, X_{q-1}) \rightarrow H_p(X_{q-1}) \rightarrow H_p(X_q) \rightarrow H_p(X_q, X_{q-1}) \rightarrow \dots .$$

The relative groups are all zero, except possibly $H_q(X_q, X_{q-1})$. Hence, the maps from $H_p(X_{q-1})$ to $H_p(X_q)$ are isomorphisms for $p+1 < q$. Composing these isomorphisms for q from $p+2$ to d implies that $H_p(X_{p+1})$ is isomorphic to $H_p(SdK)$.

The main tool in this proof is a 2-dimensional diagram connecting pieces of the long exact sequences of the pairs (X_q, X_{q-1}) for $q = p-1, p, p+1$. We write this diagram identifying $H_q(X_q, X_{q-1})$ with D_q :

$$\begin{array}{ccccc}
 D_{p+1} & & & & 0 = H_{p-1}(X_{p-2}) \\
 \downarrow e & \searrow & & & \downarrow \\
 0 = H_p(X_{p-1}) & \longrightarrow & H_p(X_p) & \xrightarrow{f} & D_p \xrightarrow{g} H_{p-1}(X_{p-1}) \\
 \downarrow l & & \searrow & & \downarrow h \\
 H_p(X_{p+1}) & & & & D_{p-1} \\
 \downarrow & & & & \\
 0 = H_p(X_{p+1}, X_p).
 \end{array}$$

We see that the long exact sequences for $q = p+1$ and $q = p-1$ run vertically, connected by the long exact sequence for $q = p$, which runs horizontally. Within this arrangement, the block chain complex runs diagonally, from the upper left to the lower right, forming two commuting diagrams. As mentioned above, the relative homology groups off the diagonal are zero, which explains the trivial group at the bottom of the diagram. We also note that $H_p(X_q) = 0$ for all $q < p$ simply because the dimension of X_q is less than p . This gives two additional trivial groups in the diagram. We are now ready for some diagram chasing, using the maps e, f, g, h, l as labeled in the diagram. The subgroup of p -cycles in D_p is the kernel of $\partial_p = h \circ g$. Since h is injective, this group is also the kernel of g . By exactness of the horizontal sequence, we have $\ker g = \text{im } f$, and since f is injective, this implies that $H_p(X_p)$ is isomorphic to the group of p -cycles. The subgroup of p -boundaries in D_p is the image of $\partial_{p+1} = f \circ e$. Since f is injective, this group is isomorphic to the image of e . The p -th homology group is the quotient of the two, $H_p(D) = H_p(X_p)/\text{im } e$. By exactness of the first vertical sequence, this is equal to $H_p(X_p)/\ker l$. But l is surjective, so this quotient is isomorphic to $H_p(X_{p+1})$ and therefore to $H_p(SdK)$, as required. \square

First form of Poincaré duality. There is a fairly direct translation between chains formed by dual blocks and cochains formed by the corresponding simplices. We have all the results lined up to prove the main result of this section.

POINCARÉ DUALITY THEOREM (FIRST FORM). Let \mathbb{M} be a compact, combinatorial d -manifold. Then there is an isomorphism between $H_p(\mathbb{M})$ and $H^q(\mathbb{M})$ for every pair of complementary dimensions $p + q = d$.

PROOF. Let K be a triangulation of \mathbb{M} with the appropriate condition on the links of its simplices, and let $p + q = d$ be two complementary dimensions. For each

p -simplex σ in K , we write σ^* for the dual p -cochain defined by $\langle \sigma^*, \sigma \rangle = 1$ and $\langle \sigma^*, \tau \rangle = 0$ for all p -simplices τ different from σ . Recall that $\hat{\sigma}$ is the dual block of σ , which is q -dimensional. To formalize the correspondence between the dual block and the dual cochain, we establish the map $\varphi_q : D_q \rightarrow C^p$ defined on the chain level by setting $\varphi_q(\hat{\sigma}) = \sigma^*$ and extending linearly. This is, of course, an isomorphism. To prove Poincaré duality, we only need to show that the introduced maps commute with the boundary and coboundary maps, that is, the diagram

$$\begin{array}{ccc} D_q & \xrightarrow{\partial_q} & D_{q-1} \\ \downarrow \varphi_q & & \downarrow \varphi_{q-1} \\ C^p & \xrightarrow{\delta^p} & C^{p+1} \end{array}$$

commutes. Going one way, we get $\varphi_{q-1} \circ \partial_q(\hat{\sigma})$, which is the $(p+1)$ -cochain that evaluates all $(p+1)$ -dimensional cofaces of σ to 1 and all other $(p+1)$ -simplices to 0. Going the other way, we get $\delta^p \circ \varphi_q(\hat{\sigma})$, which does the same. \square

Recall that the Universal Coefficient Theorem states that $H_p(M)$ is isomorphic to $H^p(M)$. Together with the Poincaré Duality Theorem, we thus have $H_p(M) \simeq H_q(M)$ for all $p + q = d$.

Bibliographic notes. Poincaré mentioned a form of his duality in a paper in 1893, without giving a proof. He tried a proof in his 1895 “Analysis situ” paper [120] based on intersection theory (see the next section), which he invented. Criticism of his work by Poul Heegard led him to realize that his proof was flawed, and he gave a new proof in two complements of the “Analysis situ” paper [121, 122], now based on dual triangulations. Poincaré duality took on its modern form in the 1930s when Eduard Čech and Hassler Whitney invented the cup and cap products of cohomology.

The proof of the Poincaré Duality Theorem presented in this section is fashioned after that of Munkres [116], simplified by the assumption that we are working with a combinatorial manifold. Not all triangulated manifolds are combinatorial, as the exotic 5-sphere described in this section shows. The proof that $\Sigma^2 P$ is homeomorphic to S^5 is due to Edwards [64]. See [125] for further exotic manifolds, including manifolds for which all triangulations violate the condition on the links. While the restriction to combinatorial manifolds is a loss of generality, the Poincaré Duality Theorem nevertheless holds for arbitrary triangulated manifolds [116]. In fact, if we use singular homology, Poincaré duality holds for arbitrary topological manifolds and even for non-compact manifolds if we use what is called cohomology with compact support. A nice proof of this can be found in [107, Chapter 20].

V.3 Intersection Theory

There is a second version of Poincarè duality, stated purely in terms of homology. It is based on an intersection pairing between homology classes of complementary dimensions introduced in this section.

Counting intersections modulo 2. Let \mathbb{M} be a combinatorial d -manifold with triangulation K . Furthermore, let p and q be integers such that $p + q = d$. As explained in Section V.2, if σ is a p -simplex in K , then its dual block, $\hat{\sigma}$, is q -dimensional. The two meet in a single point, the barycenter of σ . If τ is another p -simplex, then $\sigma \neq \tau$ implies that σ and $\hat{\tau}$ are disjoint. We therefore define

$$\sigma \cdot \hat{\tau} = \begin{cases} 1 & \text{if } \sigma = \tau; \\ 0 & \text{if } \sigma \neq \tau. \end{cases}$$

We are mainly interested in intersections of cycles. Suppose that $c = \sum_i a_i \sigma_i$ is a p -cycle in K and $d = \sum_j b_j \hat{\tau}_j$ is a q -cycle in the dual block decomposition. Then the *intersection number* of the two cycles is

$$c \cdot d = \sum_{i,j} a_i b_j (\sigma_i \cdot \hat{\tau}_j),$$

counting the intersections modulo 2. In other words, $c \cdot d = 0$ if the two cycles are disjoint or meet in an even number of points, and $c \cdot d = 1$ if they meet in an odd number of points. As an example, consider the center circle of the Möbius strip and a pulled-off copy, that is, a nearby closed curve that meets the center circle in a finite number of points, as sketched in Figure V.6. The topology of the Möbius

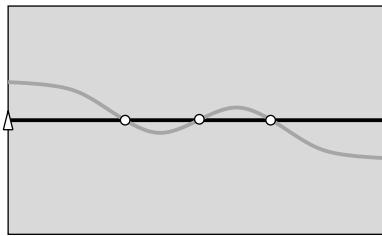


Figure V.6: The black center circle of the Möbius strip intersects the gray pulled-off copy in three points.

strip forces an odd number of intersections. This is unlike the annulus, in which a pulled-off closed curve always meets the original in an even number of points.

It is not difficult to show that if we replace c or d by a homologous cycle, then the intersection number does not change. For example, if $c \sim c_0$, we consider the intersection of d with a $(p+1)$ -chain γ in K for which $\partial\gamma = c + c_0$. Let τ be a $(p+1)$ -simplex of γ and let $\hat{\sigma}$ be a block of dimension $q = d - p$. The key observation is that τ and $\hat{\sigma}$ are disjoint unless σ is a face of τ , in which case they intersect in

the edge connecting the barycenter of τ to the barycenter of σ . Completing the intersection between γ and d , the edge extends to either a closed curve or a path with two endpoints. These points lie either both on c , or both on c_0 , or one on c and the other on c_0 . The total number of endpoints is even, which implies that the intersection numbers are the same, that is, $c \cdot d = c_0 \cdot d$.

Pairings. Since the intersection number is invariant under choosing different representatives of a homology class, we have a map $\# : H_p(\mathbb{M}) \times H_q(\mathbb{M}) \rightarrow G$ defined by $\#(\gamma, \delta) = c \cdot d$, where c and d are representative cycles of γ and δ . We call this map the *intersection pairing* of the homology groups, where $p + q = d$, as before. Using the same notation as for simplices and cycles, we write $\gamma \cdot \delta = \#(\gamma, \delta)$ and call it the *intersection number* of $\gamma \in H_p(\mathbb{M})$ and $\delta \in H_q(\mathbb{M})$. The pairing is bilinear and symmetric; that is,

$$\begin{aligned} (a\gamma + a_0\gamma_0) \cdot \delta &= a(\gamma \cdot \delta) + a_0(\gamma_0 \cdot \delta), \\ \gamma \cdot (b\delta + b_0\delta_0) &= b(\gamma \cdot \delta) + b_0(\gamma \cdot \delta_0), \\ \gamma \cdot \delta &= \delta \cdot \gamma. \end{aligned}$$

Since we work modulo 2, we do not have to worry about orientations of simplices and manifolds. To define intersection theory over an arbitrary field, we would need to deal with this issue, and the intersection number would be an element of the field. In this case, bilinearity still holds, but symmetry as stated does not.

Pairings can be defined more generally. For example, let U and V be vector spaces over $G = \mathbb{Z}_2$. A bilinear pairing $\# : U \times V \rightarrow G$ gives a natural homomorphism $\phi_{\#} : V \rightarrow \text{Hom}(U, G)$ defined by $\phi_{\#}(v) = f_v$, where $f_v(u) = u \cdot v$. The pairing is *perfect* if for every non-zero $u \in U$ there exists at least one $v_0 \in V$ with $\#(u, v_0) = 1$ and, symmetrically, for every non-zero $v \in V$ there exists at least one $u_0 \in U$ with $\#(u_0, v) = 1$.

PERFECT PAIRING LEMMA. The pairing $\# : U \times V \rightarrow G$ is perfect iff the implied natural homomorphism $\phi_{\#} : V \rightarrow \text{Hom}(U, G)$ is an isomorphism.

PROOF. Suppose first that $\phi_{\#}$ is an isomorphism. If we take $v \neq 0$, then since $\phi_{\#}$ is injective, $f_v \neq 0$, which means there is at least one u_0 with $\#(u_0, v) = 1$. Furthermore, if $u \neq 0$, since $\phi_{\#}$ is surjective, there is a $v_0 \in V$ with $\phi_{\#}(v_0) = u^*$, and this means that $\#(u, v_0) = 1$.

Conversely, suppose that the paring is perfect. The map $\phi_{\#}$ is injective because if $f_v = 0$, then $\#(u, v) = 0$ for every u , so $\#$ perfect gives $v = 0$. Note that this implies $\text{rank } V \leq \text{rank } \text{Hom}(U, G) = \text{rank } U$. The similarly defined map from U to $\text{Hom}(V, G)$ is injective by the analogous argument, which implies $\text{rank } U \leq \text{rank } \text{Hom}(V, G) = \text{rank } V$. Thus $\phi_{\#}$ is an injective map between vector spaces of the same dimension, which implies it is an isomorphism. \square

Since V and $\text{Hom}(U, G)$ are isomorphic, this implies that U and V are isomorphic. However, this isomorphism depends on a choice of basis.

Intersection and cohomology. We can define the Poincaré duality map using intersection numbers. Indeed, if σ is a p -simplex of K and $\hat{\sigma}$ is its dual block of dimension q , then $\varphi_q(\hat{\sigma}) = \sigma^*$. That is, $\varphi_q(\hat{\sigma})$ is the p -dimensional cochain for which

$$\langle \sigma^*, \tau \rangle = \begin{cases} 1 & \text{if } \sigma = \tau; \\ 0 & \text{if } \sigma \neq \tau. \end{cases}$$

Since the same holds for intersection numbers, we have $\langle \varphi_q(\hat{\sigma}), \tau \rangle = \hat{\sigma} \cdot \tau$. By linear extension, this formula holds for chains, and since it is the same for different representatives of the same class, the formula also holds for the induced map on homology, that is,

$$\langle \varphi_*(\gamma), \delta \rangle = \gamma \cdot \delta.$$

Using this formula, there is a second version of Poincaré duality.

POINCARÉ DUALITY THEOREM (SECOND FORM). Let \mathbb{M} be a compact, combinatorial d -manifold. Then the pairing $\# : H_p(\mathbb{M}) \times H_q(\mathbb{M}) \rightarrow G$ defined by $\#(\gamma, \delta) = \gamma \cdot \delta$ is perfect for all integers $p + q = d$.

The proof follows from the first form and is omitted.

The torus and the Klein bottle. To illustrate Poincaré duality formulated in terms of intersection numbers, we now consider the two examples sketched in Figure V.7. For the 2-dimensional torus, $S^1 \times S^1$, the most interesting case is in

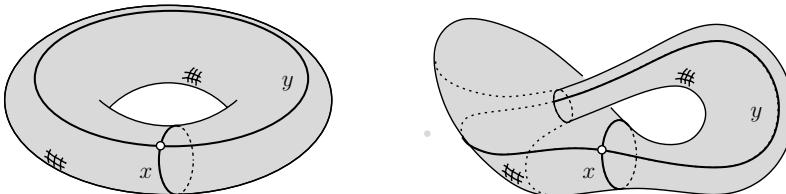


Figure V.7: The meridian and longitudinal curves of the torus on the left and of the Klein bottle on the right.

dimension 1, for which the second form of Poincaré duality gives a perfect pairing $\# : H_1 \times H_1 \rightarrow G$. Natural generators of H_1 are the *meridian curve*, x , which bounds a disk in the solid region enclosed by the torus (but not on the torus), and the *longitudinal curve*, y , which meets x in a single point and does not bound. The intersection numbers are easy to compute. Pushing off x and y gives homologous closed curves that are disjoint from the originals or meet them in an even number of points. Hence, the intersection numbers between x and x and between y and y vanish, and the intersection number between x and y is 1; see Table V.1 on the left. Note that the determinant of the matrix of intersection numbers is one.

The modulo 2 homology of the Klein bottle is the same as that of the torus. However, the intersection pairing on H_1 is different. As for the torus, we can take two curves x and y that generate H_1 with $x \cdot y = y \cdot x = 1$ and $x \cdot x = 0$. However, a neighborhood of the curve y is a Möbius strip, so pushing off y gives a closed curve that intersects y an odd number of times, that is, $y \cdot y = 1$. If we change the basis, we still do not get the same matrix as that of the torus. Once again, the matrix of intersection numbers, given in Table V.1 on the right, has determinant one.

	x	y	x	y
x	0	1	0	1
y	1	0	1	1

Table V.1: The intersection numbers of the meridian and the longitudinal curves for the torus on the left and the Klein bottle on the right.

Euler characteristic. By the Euler-Poincaré Theorem, the Euler characteristic of any space is the alternating sum of its Betti numbers. Letting \mathbb{M} be a compact, combinatorial d -manifold, the Poincaré Duality and the Universal Coefficient Theorems imply $\beta_i = \beta_{d-i}$ for all i . For odd d , this gives

$$\chi(\mathbb{M}) = \beta_0 - \beta_1 + \dots + \beta_{d-1} - \beta_d,$$

which vanishes. For even d , this tells us that the terms above and below half the dimension contribute equal amounts to the Euler characteristic. Writing $d = 2k$, this gives

$$\chi(\mathbb{M}) = 2[\beta_0 - \beta_1 + \dots \pm \beta_{k-1}] \mp \beta_k.$$

It follows that the Euler characteristic is even iff β_k is even. However, the group $H_k(\mathbb{M})$ is paired with itself and is therefore self-dual. If \mathbb{M} is orientable, this can be used to show that β_k and the Euler characteristic are both even. In contrast, homology and cohomology modulo 2 do not capture this subtlety.

Manifolds with boundary. If \mathbb{M} is a manifold with boundary, Poincaré duality does not hold. For example, if we take the ball, \mathbb{B}^d , its 0-dimensional homology has rank one while its d -dimensional homology vanishes. There is a form for manifolds with boundary, however, called Lefschetz duality, which reduces to Poincaré duality when the boundary is empty. It relates an absolute homology or cohomology group to a relative one. Returning to our example, note that $H_0(\mathbb{B}^d)$ and $H_d(\mathbb{B}^d, \mathbb{S}^{d-1})$ both have rank one.

LEFSCHETZ DUALITY THEOREM (FIRST FORM). Let \mathbb{M} be a compact, combinatorial d -manifold with boundary $\partial\mathbb{M}$. Then for every pair of complementary dimensions, $p + q = d$, there are isomorphisms $H_p(\mathbb{M}, \partial\mathbb{M}) \simeq H^q(\mathbb{M})$ and $H_p(\mathbb{M}) \simeq H^q(\mathbb{M}, \partial\mathbb{M})$.

Again this can be combined with the Universal Coefficient Theorem, $H_p(M) \simeq H^p(M)$, to see that $H_p(M, \partial M) \simeq H_q(M)$ for all $p + q = d$. The proof of the Lefschetz Duality Theorem follows that of the Poincaré Duality Theorem exactly, inserting relative chains and cochains where needed, so we omit the proof. There is also a second version of Lefschetz duality based on the extension of the intersection pairing to a pairing between absolute and relative classes. Again we omit the details and the proof.

LEFSCHETZ DUALITY THEOREM (SECOND FORM). Let M be a compact, combinatorial d -manifold with boundary ∂M . Then the intersection pairing $\# : H_p(M) \times H_q(M, \partial M) \rightarrow G$ is perfect for all $p + q = d$.

We illustrate Lefschetz duality formulated in terms of intersection numbers for the half torus sketched in Figure V.8. Being homeomorphic to the cylinder, the first homology group has a single generator, the meridian curve of the full torus. Similarly, the first relative homology group has a single generator, namely half the longitudinal curve; see Figure V.8.

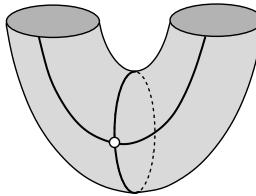


Figure V.8: The displayed generators of the first absolute and first relative homology groups of the half torus meet in a single point.

Bibliographic Notes. Henri Poincaré invented intersection theory to prove his duality theorem in 1895 [120], but this attempt failed. It is also said that Alexander and Lefschetz founded the intersection theory of cycles on manifolds in the 1920s. Their theory was one of the precursors of cohomology. The Lefschetz Duality Theorem dates back to the 1920s when Solomon Lefschetz introduced it along with the concept of relative homology [101]. The statements of the theorem given in this section are limited to combinatorial manifolds. The sole reason is our desire to keep the proof simple. Indeed, the theorem holds in more generality for manifolds with boundary. A good modern account of the theorem can be found in [106].

V.4 Alexander Duality

Prisons in d -dimensional space are made of $(d-1)$ -dimensional walls. This is because a wall of lower dimension cannot separate space. The topic of this section is a formal expression of a generalization of this statement and its use in the design of a fast algorithm for homology.

Separating water from land. The 2-sphere decomposes the 3-sphere into two balls, both homologically trivial. Compare this with the torus decomposing the 3-sphere into two solid tori, both with non-trivial first homology. The two examples suggest a relationship between the homology of a subspace and its complement. Such a relationship exists for any manifold, but the prettiest case is when the manifold is a sphere. Let, therefore, K be a triangulation of \mathbb{S}^d . To simplify the technical discussions, we assume that the link of every simplex is a sphere of the appropriate dimension. Let B be the dual block decomposition. We call $N \subseteq K$ and $X \subseteq B$ *complementary subcomplexes* if a simplex $\sigma \in K$ belongs to N iff its dual block, $\hat{\sigma} \in B$, does not belong to X ; see Figure V.9. Recall that the dual blocks are

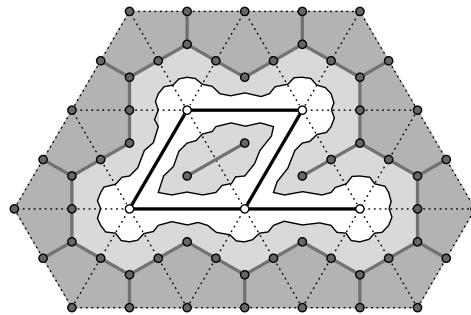


Figure V.9: A portion of a triangulation of the 2-sphere. We see a subcomplex, N , consisting of five edges and five vertices. It is surrounded by the complementary subcomplex of the dual block decomposition, X . The two are separated by two curves in the 1-skeleton of the second barycentric subdivision.

defined as subcomplexes of the first barycentric subdivision, SdK . To separate N from X , we subdivide once more, and we note that Sd^2K contains subcomplexes N' and X' whose underlying spaces are the same as those of N and X . We enlarge both by adding a layer of simplices, making sure we still have complexes. Specifically, let

$$\begin{aligned} N'' &= \bigcup_{u \in N'} \overline{\text{St}} u, \\ X'' &= \bigcup_{v \in X'} \overline{\text{St}} v; \end{aligned}$$

see Figure V.9. Since each vertex in the first barycentric subdivision belongs to either N' or to X' , the complexes N'' and X'' exhaust the second barycentric subdivision. Indeed, $|N''|$ and $|X''|$ are d -manifolds with disjoint interiors and common boundary, $\partial N'' = \partial X'' = N'' \cap X''$, which is a $(d-1)$ -manifold separating N and X . This motivates us to call X'' the *exterior* of the decomposition. Note also that there are deformation retractions from $|N''|$ to $|N|$ and from $|X''|$ to $|X|$.

The generalized prison wall theorem. Traditionally, the mentioned relationship is stated in terms of the homology and cohomology groups of spaces $\mathbb{N} \subseteq \mathbb{S}^d$ and

$\mathbb{X} = \mathbb{S}^d - \mathbb{N}$. If \mathbb{N} is closed, then \mathbb{X} is open. Since we did not define homology groups for open sets, we state the relationship in terms of complementary subcomplexes.

ALEXANDER DUALITY THEOREM. Let \mathbb{S}^d be a combinatorial d -sphere with triangulation K , N a subcomplex of K , and X the complementary subcomplex of the dual block decomposition. Then $\tilde{H}_p(N) \simeq \tilde{H}^{d-p-1}(X)$ for all dimensions p .

PROOF. We first consider the case $p < d-1$, by showing a sequence of isomorphisms:

$$\tilde{H}^{d-p-1}(X) \simeq \tilde{H}^{d-p-1}(X'') \quad (\text{V.1})$$

$$\simeq H^{d-p-1}(X'') \quad (\text{V.2})$$

$$\simeq H_{p+1}(X'', \partial X'') \quad (\text{V.3})$$

$$\simeq \tilde{H}_{p+1}(Sd^2K, N'') \quad (\text{V.4})$$

$$\simeq \tilde{H}_p(N'') \quad (\text{V.5})$$

$$\simeq \tilde{H}_p(N). \quad (\text{V.6})$$

We get equation (V.1) because $|X|$ is a deformation retract of $|X''|$, equation (V.2) because cohomology and reduced cohomology are the same in dimension $d-p-1 > 0$, equation (V.3) by Lefschetz duality for X'' , which is a d -manifold with boundary $\partial X''$, and equation (V.4) by excision. We get equation (V.5) by considering the long exact sequence of the pair for reduced homology,

$$\dots \rightarrow \tilde{H}_{p+1}(Sd^2K) \rightarrow H_{p+1}(Sd^2K, N'') \rightarrow \tilde{H}_p(N'') \rightarrow \tilde{H}_p(Sd^2K) \rightarrow \dots,$$

and noticing that the p -th and $(p+1)$ -st homology groups of the sphere are trivial because $p+1 < d$. Finally, we get equation (V.6) because $|N|$ is a deformation retract of $|N''|$.

The case $p = d-1$ is similar. We again have the sequence of six isomorphisms and similar reasons for each. In equation (V.1), we have an extra copy of G on both sides, which we lose in equation (V.2) because of the difference between ordinary and reduced cohomology in dimension zero. We pick up the copy again in equation (V.5) because the rank of $\tilde{H}_d(Sd^2K)$ in the long exact sequence is one. This same copy of G also appears in equation (V.6). Finally, for $p = -1$ or d , we have either two trivial groups, which are therefore isomorphic, or two groups both isomorphic to G , namely when $N = \emptyset$ or $B = \emptyset$. \square

The sole reason for limiting the theorem to the combinatorial d -sphere is the use of the Lefschetz Duality Theorem in its proof. Since the latter holds for general triangulations, so does the Alexander Duality Theorem. There are many applications of this theorem, including the application to knots in \mathbb{R}^3 . We focus on the computation of Betti numbers.

Adding a simplex. We begin by studying how the addition of a single simplex affects the ranks of the homology groups. Let $N_{i-1} \subseteq N_i$ be simplicial complexes

that differ by a single simplex, that is, $N_i - N_{i-1} = \{\sigma_i\}$. Consider the long exact reduced homology sequence of the pair,

$$\dots \rightarrow \tilde{H}_p(N_{i-1}) \xrightarrow{\varphi} \tilde{H}_p(N_i) \xrightarrow{D} H_p(N_i, N_{i-1}) \rightarrow \tilde{H}_{p-1}(N_{i-1}) \xrightarrow{\psi} \tilde{H}_{p-1}(N_i) \rightarrow \dots,$$

where D is the connecting homomorphism in dimension $p = \dim \sigma_i$. The group $H_q(N_i, N_{i-1})$ is trivial for all $q \neq p$ and has rank one for $q = p$. It follows that the maps from the reduced homology groups of N_{i-1} to those of N_i are isomorphisms, except for $\varphi : \tilde{H}_p(N_{i-1}) \rightarrow \tilde{H}_p(N_i)$, which is possibly only injective, and $\psi : \tilde{H}_{p-1}(N_{i-1}) \rightarrow \tilde{H}_{p-1}(N_i)$, which is possibly only surjective. There are two possibilities for D .

CASE 1: D is surjective. Then $\tilde{\beta}_p(A) = \tilde{\beta}_p(A_0) + 1$.

CASE 2: D is the zero map. Then $\tilde{\beta}_{p-1}(A) = \tilde{\beta}_{p-1}(A_0) - 1$.

All other Betti numbers remain the same. In Case 1, the addition of σ creates a new homology class, so we call it a *positive simplex*. In Case 2, σ destroys a homology class, so we call it a *negative simplex*. The difference between the two cases will be important again later, when we define persistent homology in Chapter VII.

Incremental algorithm. We use the insight about the two types of simplices to compute the Betti numbers of a simplicial complex by adding one simplex at a time. Let $\sigma_1, \sigma_2, \dots, \sigma_j$ be an ordering of the simplices in N such that every prefix of the sequence forms a complex. In other words, $N_i = \{\sigma_1, \sigma_2, \dots, \sigma_i\}$ is a subcomplex of N for $0 \leq i \leq j$. The algorithm starts with $N_0 = \emptyset$, whose reduced Betti numbers are zero except in dimension minus one.

```

 $\tilde{\beta}_{-1} = 1;$  for  $p = 0$  to  $d$  do  $\tilde{\beta}_p = 0$  endfor;
for  $i = 1$  to  $j$  do
  if  $\sigma_i$  is positive then  $\tilde{\beta}_p = \tilde{\beta}_p + 1$ 
  else  $\tilde{\beta}_{p-1} = \tilde{\beta}_{p-1} - 1$ 
  endif
endfor.

```

We call this the Incremental Betti Number Algorithm. Assuming we know the classification of the simplices, the algorithm computes the Betti numbers of all complexes N_i spending only constant time per simplex. We emphasize that the algorithm makes no assumption on the dimension or the topology of the complex, other than that both are finite. The only difficult operation is the classification of the simplices. We could adapt the matrix reduction algorithm for homology described in Section IV.2, and we will revisit this idea in Chapter VII. However, in some important cases, the classification can be done directly, without the use of the boundary matrix.

Fast classification in \mathbb{S}^3 . Suppose K is a triangulation of \mathbb{S}^3 and $\sigma_1, \sigma_2, \dots, \sigma_m$ is an ordering of its simplices such that $N_i = \{\sigma_1, \sigma_2, \dots, \sigma_i\}$ is a subcomplex of

K for $0 \leq i \leq m$. By stereographic projection, this includes the case in which K is the Delaunay complex of a finite set of points in \mathbb{R}^3 and the N_i are the alpha complexes of the same set; see Chapter III. We classify the simplices in the order of their dimension, reversing two and three to exploit duality.

Vertices. The first vertex of K is negative because adding it reduces $\tilde{\beta}_{-1}$ by one. All other vertices are positive.

Edges. An edge is negative if it merges two components into one, and it is positive if it connects two vertices that already belong to a common component; see Figure V.10. To distinguish between the two cases, we maintain the components in a union-find data structure, as described in Section I.1. If σ_i is an edge, we label it negative if its two endpoints lie in different components of N_{i-1} , in which case the addition of σ_i triggers a union operation. Otherwise, we label σ_i positive.

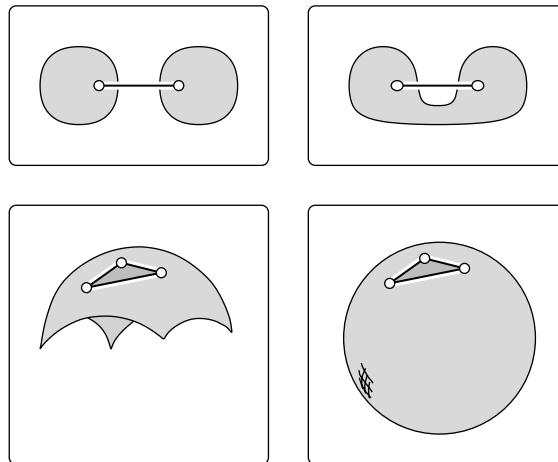


Figure V.10: Top: a negative edge on the left and a positive edge on the right. Bottom: a negative triangle on the left and a positive triangle on the right.

We note that the classification of the vertices and edges as described works in general, for any finite simplicial complex, K . In contrast, the classification of triangles and tetrahedra below makes crucial use of the assumption that K triangulates the 3-sphere.

Tetrahedra. The last tetrahedron of K is positive, and all other tetrahedra are negative.

Triangles. A triangle is negative if it closes a tunnel and positive if it forms a void; see Figure V.10. This is easier to explain for the dual blocks. Let X_i be the set of dual blocks of the simplices in $K - N_i$ and note that it is a subcomplex of B . Reading the simplices from back to front and replacing them by their dual blocks is like enumerating the sets from $X_m = \emptyset$ to $X_0 = K$. As before,

we maintain the components in a union-find data structure. Assuming $\hat{\sigma}_i$ is an edge, we label σ_i negative if the two endpoints of $\hat{\sigma}_i$ lie in different components of X_{i+1} , in which case the addition of $\hat{\sigma}_i$ triggers a union operation. Otherwise, we label σ_i positive.

It is straightforward to see that the classification of the vertices and edges is correct. For the tetrahedra, we use the fact that K triangulates \mathbb{S}^3 . For the triangles, we use Alexander duality. As discussed in Section I.1, the maintenance of a union-find data structure takes time proportional to $\alpha(n)$ per operation, where n is the number of nodes (the vertices or tetrahedra in K), and α is the inverse of the Ackermann function. For all practical purposes, this is constant time per operation.

Bibliographic note. We may think of the Alexander Duality Theorem as a generalization of the Jordan Curve Theorem mentioned in Chapter I. Roots of the concept can be found in the work of Alexander [8], which was later further developed, in particular by Alexandrov and Pontryagin. Similar to the Poincaré and the Lefschetz Duality Theorems, we state and prove the Alexander Duality Theorem only for the combinatorial manifold case. The restriction is unnecessary but simplifies the exposition.

The Incremental Betti Number Algorithm is due to Delfinado and Edelsbrunner [45]. It was implemented as part of the 3-dimensional Alpha Shape software written and distributed by Ernst Mücke in the early nineties of the last century. The algorithm exploits the structure of the alpha shape filtration [63] and introduces the concept of positive and negative simplices, thus foreshadowing the later development of persistent homology, which will be discussed in Chapter VII.

Exercises

The credit assignment reflects a subjective assessment of difficulty. A typical question can be answered using knowledge of the material combined with some thought and analysis.

1. **Coboundary** (one credit). Prove that the coboundary map can be thought of as taking each simplex to its cofaces of one dimension higher. Formally, $\langle \delta\varphi, \tau \rangle = 1$ iff $\langle \varphi, \sigma \rangle = 1$ for an odd number of faces σ of τ with dimension $\dim \sigma = \dim \tau - 1$.
2. **Universal Coefficient Theorem** (two credits). Let $\varphi \in Z^p$ be a cocycle representing a cohomology class $\gamma \in H^p$ and let $c \in Z_p$ be a cycle representing a homology class $\alpha \in H_p$. Let $j : H^p \rightarrow \text{Hom}(H_p, \mathbb{Z}_2)$ be defined so that $j(\gamma)$ applied to α is equal to $\langle \varphi, c \rangle$.
 - (i) Show that j is well defined, that is, it does not depend on the representatives chosen for γ and α .
 - (ii) Show that j is an isomorphism.

3. **Dual vector spaces** (three credits). Let U be a vector space over $G = \mathbb{Z}_2$ and let $U^* = \text{Hom}(U, G)$ be its dual.

- (i) Show that U^* is also a vector space and U and U^* are isomorphic. However, note that the isomorphism between U and U^* depends on a choice of basis and is thus not natural.

Let $(U^*)^* = \text{Hom}(U^*, G)$ be the dual of the dual of U . Let $j : U \rightarrow (U^*)^*$ be defined by mapping $u \in U$ to $j(u) = \phi \in (U^*)^*$ such that $\phi(f) = f(u)$ for every $f \in U^*$.

- (ii) Prove that j is an isomorphism.

4. **Classifying 3-manifolds** (one credit). Recall the two steps of the algorithm for classifying a triangulated 2-manifold described in Chapter II: deciding orientability and computing the Euler characteristic.

- (i) Generalize both steps of the algorithm to a triangulated 3-manifold.
(ii) Why does the algorithm in (i) not suffice to classify 3-manifolds?

5. **Poincaré duality** (two credits). Use the Perfect Pairing Lemma to prove the first form from the second form of the Poincaré Duality Theorem.

6. **2-manifold with boundary** (one credit). Let M be a 2-manifold with genus g and b boundary circles.

- (i) Use the first form of the Lefschetz Duality Theorem to determine the ranks of $H_p(M)$ and $H_p(M, \partial M)$ for $p = 0, 1, 2$.
(ii) Draw generators for the first absolute and relative homology groups and check your drawing against the second form of the Lefschetz Duality Theorem.

7. **d -dimensional boundary** (three credits). Let M be a $(d+1)$ -manifold and ∂M its boundary, a d -manifold without boundary.

- (i) Use Lefschetz and Poincaré duality to show that the kernel of the map $f_d : H_d(\partial M) \rightarrow H_d(M)$ induced by the inclusion $\partial M \subseteq M$ has rank exactly half the d -th Betti number of ∂M .
(ii) Show that $x \cdot y = 0$ for all classes x and y in $\ker f_d$.

8. **Water and land on a manifold** (one credit). Let M be a combinatorial d -manifold and K its triangulation. Let N be a subcomplex of K , N' the subcomplex of SdK whose underlying space is the same as that of N , and N'' the smallest subcomplex of Sd^2K that contains all simplices incident to vertices in N' .

- (i) Show that N'' is the union of closed stars over all vertices of N' .
(ii) Show that N'' is a d -manifold with boundary.

Chapter VI

Morse Functions

The class of real-valued functions on a manifold is an unwieldy animal, and restricting it to continuous functions does not do a whole lot to tame it. Even smooth functions can be rather complicated in their behavior, and it is best to add another requirement, namely genericity. What we get then is the class of Morse functions, which distinguishes itself by having only simple critical points. Most of the theory is concerned with the study of these critical points, their structure, and what they say about the manifold and the function. In spite of the fact that we rarely find Morse functions in actual applications, or smooth functions for that matter, knowing about their structure significantly benefits our understanding of general, smooth functions and even piecewise linear functions.

VI.1 Generic Smooth Functions

Many questions in the sciences and engineering are posed in terms of real-valued functions. Instead of struggling with the wild character of general functions, we restrict our attention to a class that achieves structural simplicity without undue limitation of shape.

The upright torus. We start with an example that foreshadows many of the results on generic smooth functions. Let \mathbb{M} be the 2-dimensional torus and $f(x)$ the height of the point $x \in \mathbb{M}$ above a horizontal plane on which the torus rests, as in Figure VI.1. We call $f : \mathbb{M} \rightarrow \mathbb{R}$ a *height function*. Each real number a has a preimage, $f^{-1}(a)$, which we refer to as a *level set*. It consists of all points $x \in \mathbb{M}$ at height a . Accordingly, the *sublevel set* consists of all points at height at most a :

$$\mathbb{M}_a = f^{-1}(-\infty, a] = \{x \in \mathbb{M} \mid f(x) \leq a\}.$$

We are interested in the evolution of the sublevel set as we increase the level a . Critical events occur when a passes the height values of the points u, v, w, z in

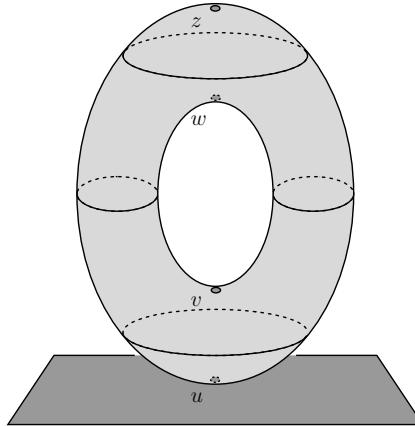


Figure VI.1: The vertical height function on the torus with critical points u , v , w , z and level sets between their height values.

Figure VI.1. For $a < f(u)$, the sublevel set is empty. For $f(u) < a < f(v)$, it is a disk, which has the homotopy type of a point. For $f(v) < a < f(w)$, the sublevel set is a cylinder. It has the homotopy type of a circle. For $f(w) < a < f(z)$, the sublevel set is a torus with a disk removed. It has the homotopy type of a figure-eight curve. Finally, for $f(z) < a$, we have the complete torus. It is obtained by gluing on the final disk. Figure VI.2 illustrates the three intermediate stages of the evolution. We need some background in differential topology to explain in what sense this evolution of the sublevel set is representative of the general situation.

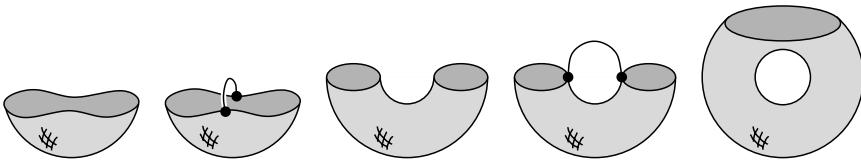


Figure VI.2: Going from a disk to a cylinder is homotopically the same as attaching a 1-cell. Similarly, going from the cylinder to the capped torus is homotopically the same as attaching another 1-cell.

Smooth functions. Let \mathbb{M} be a smooth d -manifold; that is, \mathbb{M} has an atlas of coordinate charts each diffeomorphic to an open ball in \mathbb{R}^d . We recall that a diffeomorphism is a homeomorphism that is smooth in both directions. Technically, being smooth means that derivatives of all orders exist. Practically, we just need derivatives of first and second order for most of the things we do, but it is easier to assume than to keep books. Denote the tangent space at a point $x \in \mathbb{M}$ by $T\mathbb{M}_x$. It is the d -dimensional vector space consisting of all tangent vectors of \mathbb{M} at x . A smooth mapping to another smooth manifold, $f : \mathbb{M} \rightarrow \mathbb{N}$, induces a linear

mapping between the tangent spaces, the derivative $Df_x : TM_x \rightarrow TN_{f(x)}$. We are primarily interested in real-valued functions for which $\mathbb{N} = \mathbb{R}$. Accordingly, we have linear maps $Df_x : TM_x \rightarrow T\mathbb{R}_{f(x)}$. The tangent space at a point of the real line is again a real line, so this is just a fancy way of saying that the derivatives are real-valued linear maps on the tangent spaces. Being linear, the image of such a map is either the entire line or just zero. We call $x \in M$ a *regular point* of f if Df_x is surjective and we call x a *critical point* of f if Df_x is the zero map. If we have a local coordinate system (x_1, x_2, \dots, x_d) in a neighborhood of x , then x is critical iff all its first-order partial derivatives vanish:

$$\frac{\partial f}{\partial x_1}(x) = \frac{\partial f}{\partial x_2}(x) = \dots = \frac{\partial f}{\partial x_d}(x) = 0.$$

The chain rule tells us that whether a point is critical or not is independent of coordinates. The image of a critical point, $f(x)$, is called a *critical value* of f . All others are *regular values* of f . We use second derivatives to further distinguish between different types of critical points. The *Hessian* of f at the point x is the matrix of second derivatives:

$$H(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1}(x) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(x) & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_d}(x) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(x) & \frac{\partial^2 f}{\partial x_2 \partial x_2}(x) & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_d}(x) \\ \vdots & \ddots & & \vdots \\ \frac{\partial^2 f}{\partial x_d \partial x_1}(x) & \frac{\partial^2 f}{\partial x_d \partial x_2}(x) & \dots & \frac{\partial^2 f}{\partial x_d \partial x_d}(x) \end{bmatrix}.$$

A critical point x is *non-degenerate* if the Hessian is non-singular, that is, $\det H(x) \neq 0$. Again, this does not depend on coordinates. The points u, v, w, z in Figure VI.1 are examples of non-degenerate critical points. Examples of degenerate critical points are $x_1 = 0$ for the function $f : \mathbb{R} \rightarrow \mathbb{R}$ defined by $f(x_1) = x_1^3$ and $(x_1, x_2) = (0, 0)$ for $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by $f(x_1, x_2) = x_1^3 - 3x_1x_2^2$. The degenerate critical point in the latter example is often referred to as a monkey saddle. Indeed, the graph of the function in a neighborhood goes up and down three times, providing a convenient resting place for the two legs as well as the tail of the monkey.

Morse functions. At a critical point, all first-order partial derivatives vanish. A local Taylor expansion has therefore no linear terms. If the critical point is non-degenerate, then the behavior of the function in a small neighborhood is dominated by the quadratic terms. Furthermore, we can find local coordinates such that there are no higher-order terms.

MORSE LEMMA. Let u be a non-degenerate critical point of $f : M \rightarrow \mathbb{R}$. There are local coordinates with $u = (0, 0, \dots, 0)$ such that

$$f(x) = f(u) - x_1^2 - \dots - x_q^2 + x_{q+1}^2 + \dots + x_d^2$$

for every point $x = (x_1, x_2, \dots, x_d)$ in a small neighborhood of u .

The number of minus signs in the quadratic polynomial is independent of coordinates and is called the *index* of the critical point, $\text{index}(u) = q$. The index classifies the non-degenerate critical points into $d + 1$ types. For a 2-manifold, we have three types, *minima* with index 0, *saddles* with index 1, and *maxima* with index 2. Examples of all three types can be seen in Figure VI.1. In Figure VI.3, we display them by showing the local evolution of the sublevel set. A consequence of the

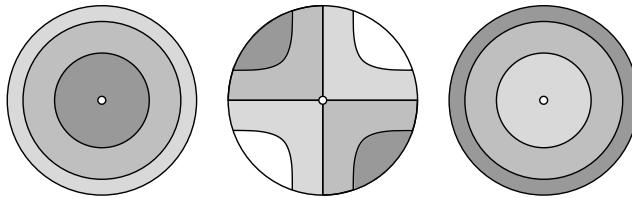


Figure VI.3: From left to right: the local pictures of a minimum, a saddle, and a maximum. Imagine looking from above with the shading getting darker as the function shrinks away from the viewpoint.

Morse Lemma is that non-degenerate critical points are isolated. In other words, each critical point has a local neighborhood that separates it from the others. This implies that a Morse function on a compact manifold has at most a finite number of critical points. To contrast this with a function that is not Morse, take the height function of a torus, similar to Figure VI.1 but placing the torus sideways, the way it would naturally rest under the influence of gravity. This height function has an entire circle of minima and another circle of maxima. All these critical points are degenerate, and their index is not defined.

DEFINITION. A *Morse function* is a smooth function on a manifold, $f : \mathbb{M} \rightarrow \mathbb{R}$, such that (i) all critical points are non-degenerate and (ii) the critical points have distinct function values.

Sometimes the second condition is dropped, but in this book we will always require both. For a geometrically perfect torus, the height function satisfies condition (i) for all but two directions, the ones parallel to the symmetry axis of the torus. Condition (ii) is violated for another two circles of directions along which the two saddles have the same height. The height function of \mathbb{S}^2 is a Morse function for all directions. The distance from a point is a Morse function for almost all points. Exceptions for the torus are points on the symmetry axis and on the center circle, but there are others. The only exception for the 2-sphere is the center.

Gradient vector field. A *vector field* on a manifold is a function $X : \mathbb{M} \rightarrow T\mathbb{M}$ that maps every point $x \in \mathbb{M}$ to a vector $X(x)$ in the tangent space of \mathbb{M} at x . Given $f : \mathbb{M} \rightarrow \mathbb{R}$ and X , we denote the directional derivative of f along the vector field by $X[f]$. It maps every point $x \in \mathbb{M}$ to the derivative of f at x in the direction $X(x)$. A particularly useful vector field is the one that points in the direction of steepest increase. To define it, we need to measure length, which we

do by introducing a Riemannian metric, that is, a smoothly varying inner product defined on the tangent spaces. For example, if \mathbb{M} is smoothly embedded in some Euclidean space, then the tangent spaces are linear subspaces of the same Euclidean space and we can borrow the metric. Given a smooth manifold \mathbb{M} , a Riemannian metric on \mathbb{M} , and a smooth function $f : \mathbb{M} \rightarrow \mathbb{R}$, we define the *gradient* of f as the vector field $\nabla f : \mathbb{M} \rightarrow T\mathbb{M}$ characterized by $\langle X(x), \nabla f(x) \rangle = X[f]$ for every vector field X . Assuming local coordinates with orthonormal unit vectors x_i , the gradient at the point x is

$$\nabla f(x) = \left[\frac{\partial f}{\partial x_1}(x), \frac{\partial f}{\partial x_2}(x), \dots, \frac{\partial f}{\partial x_d}(x) \right]^T.$$

We use the gradient to introduce a *1-parameter group of diffeomorphisms* $\varphi : \mathbb{R} \times \mathbb{M} \rightarrow \mathbb{M}$. There are two characteristic properties of this group. First, the map $\varphi_t : \mathbb{M} \rightarrow \mathbb{M}$ defined by $\varphi_t(x) = \varphi(t, x)$ is a diffeomorphism of \mathbb{M} to itself for each $t \in \mathbb{R}$, and second, $\varphi_{t+t_0} = \varphi_t \circ \varphi_{t_0}$ for all $t, t_0 \in \mathbb{R}$. Such a group defines a vector field by differentiation, and we require that this vector field be parallel to the gradient vector field:

$$\lim_{\varepsilon \rightarrow 0} \frac{f(\varphi_\varepsilon(x)) - f(x)}{\varepsilon} = \frac{\nabla f(x)}{\|\nabla f(x)\|^2}[f].$$

This group of diffeomorphisms follows the evolution of the sublevel set and can be used to prove that there are no topological changes that happen between contiguous critical values. Specifically, let $f : \mathbb{M} \rightarrow \mathbb{R}$ be smooth and let $a < b$ be such that $f^{-1}[a, b]$ is compact and contains no critical points of f . Then \mathbb{M}_a is diffeomorphic to \mathbb{M}_b .

Attaching cells. The situation is different when we consider regular values $a < b$ such that $f^{-1}[a, b]$ is compact but contains one critical point of f . Let this critical point be u and let its index be q . In this case, \mathbb{M}_b has the homotopy type of \mathbb{M}_a with a q -cell attached along its boundary. To explain what this means, we recall that \mathbb{B}^q is the q -dimensional unit ball with \mathbb{S}^{q-1} as its boundary. Let $g : \mathbb{S}^{q-1} \rightarrow \text{bd } \mathbb{M}_a$ be a continuous map. To *attach* the cell to \mathbb{M}_a , we identify each point $x \in \mathbb{S}^{q-1}$ with its image $g(x) \in \text{bd } \mathbb{M}_a$. The only case that is a bit different is $q = 0$. Then \mathbb{S}^{-1} is empty and attaching the 0-cell just means adding a point.

We illustrate this construction for a 3-manifold, \mathbb{M} . There are four types of critical points, namely minima with index 0, saddles with index 1 or 2, and maxima with index 3. The two types of saddles deserve some attention. To illustrate the local evolution of the sublevel set, we draw spheres around them and shade the portion that belongs to the sublevel set, as in Figure VI.4. The level set that passes through a saddle forms locally a double-cone with the apex at the critical point. This is the same for both types, the only difference being the side on which the sublevel set resides. For the index-1 saddle, we imagine a two-sheeted hyperboloid approaching from two sides until the two sheets meet at the saddle. Thereafter, the sublevel set thickens around the saddle as its boundary continues as a one-sheeted hyperboloid (an hour glass). Homotopically, this evolution is the same as attaching

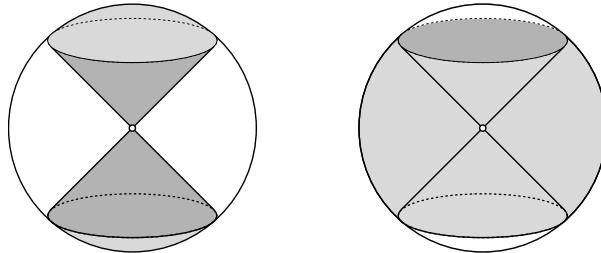


Figure VI.4: The double-cone neighborhood of the index-1 saddle on the left and of the index-2 saddle on the right. The volume occupied by the sublevel set is shaded.

a 1-cell connecting the two sheets. For the index-2 saddle, the sequence of events is reversed. Specifically, a one-sheeted hyperboloid approaches along a circle of directions until it reaches the saddle. Thereafter, the sublevel set thickens around the saddle as its boundary continues as two sheets of a hyperboloid. Homotopically, this evolution is the same as attaching a 2-cell closing the tunnel formed by the one-sheeted hyperboloid.

Bibliographic notes. Morse theory developed first in infinite dimensions, as part of the calculus of variations; see Morse [113]. The classic source on the subject for finite-dimensional manifolds is the text by Milnor [111], but see also Matsumoto [105] and Banyaga and Hurtubis [14].

VI.2 Transversality

Given a Morse function, we can follow the gradient flow and decompose the manifold depending on where the flow originates and where it ends. For this decomposition to form a complex, we require that the function satisfy an additional genericity assumption.

Integral lines. Recall the 1-parameter group of diffeomorphisms, $\varphi : \mathbb{R} \times M \rightarrow M$, defined by a Morse function, f , on a compact manifold, M , with a Riemannian metric. The *integral line* that passes through a regular point, $x \in M$, is $\gamma = \gamma_x : \mathbb{R} \rightarrow M$ defined by $\gamma(t) = \varphi(t, x)$; see Figure VI.5. It is the solution to the ordinary differential equation defined by $\dot{\gamma}(t) = \nabla f(\gamma(t))$ and the initial condition $\gamma(0) = x$. Because φ and therefore γ are defined for all $t \in \mathbb{R}$ and M is compact, the integral line necessarily approaches a critical point, both for t going to plus and to minus infinity. We call these critical points the *origin* and the *destination* of the integral line:

$$\text{org}(\gamma) = \lim_{t \rightarrow -\infty} \gamma(t), \quad \text{dest}(\gamma) = \lim_{t \rightarrow \infty} \gamma(t).$$

The function increases along the integral line, which implies that $\text{org}(\gamma) \neq \text{dest}(\gamma)$. The Existence and Uniqueness Theorems of ordinary differential equations imply that the integral line that passes through another regular point y is either disjoint from or the same as the one passing through x , $\text{im } \gamma_x = \text{im } \gamma_y$ or $\text{im } \gamma_x \cap \text{im } \gamma_y = \emptyset$. This property suggests we decompose the manifold into integral lines or unions of integral lines with shared characteristics.

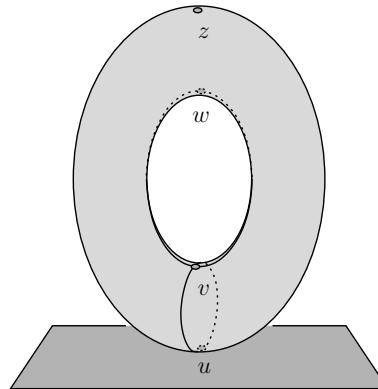


Figure VI.5: The upright torus with the four integral lines that end at the two saddles.

Stable and unstable manifolds. The *stable manifold* of a critical point u of f is the point itself together with all regular points whose integral lines end at u . Symmetrically, the *unstable manifold* of u is the point itself together with all regular points whose integral lines originate at u . More formally,

$$\begin{aligned} S(u) &= \{u\} \cup \{x \in \mathbb{M} \mid \text{dest}(\gamma_x) = u\}, \\ U(u) &= \{u\} \cup \{y \in \mathbb{M} \mid \text{org}(\gamma_y) = u\}. \end{aligned}$$

The function increases along integral lines. It follows that $f(u) \geq f(x)$ for all points x in the stable manifold of u . This is the reason why $S(u)$ is sometimes referred to as the *descending manifold* of u . Symmetrically, $f(u) \leq f(y)$ for all points y in the unstable manifold of u , and $U(u)$ is sometimes referred to as the *ascending manifold* of u .

Suppose the dimension of \mathbb{M} is d and the index of the critical point u is q . Then there is a $(q - 1)$ -sphere of directions along which integral lines approach u . It can be proved that together with u , these integral lines form an open ball of dimension q and that $S(u)$ is a submanifold homeomorphic to \mathbb{R}^q that is immersed in \mathbb{M} . It is not embedded because distant points in \mathbb{R}^q may map to arbitrarily close points in \mathbb{M} , as we can see in Figure VI.5. For example, the saddle v has a stable 1-manifold consisting of two integral lines that merge at v to form one open, connected interval. The two ends of the interval approach the minimum, u , which does not belong to the 1-manifold. While the map from \mathbb{R}^1 to \mathbb{M} is continuous, the inverse defined on its image is not.

Morse-Smale functions. The stable manifolds do not necessarily form a complex. Specifically, it is possible that the boundary of a stable manifold is not the union of other stable manifolds of lower dimension. Take for example the upright torus in Figure VI.5. The stable 1-manifold of the upper saddle, w , reaches down to the lower saddle, v , but the latter is not a stable 0-manifold. The reason for this deficiency is a degeneracy in the gradient flow. In particular, we have an integral line that originates at a saddle and ends at another saddle. Equivalently, the integral line belongs to the stable 1-manifold of w and to the unstable 1-manifold of v . Generically, such integral lines do not exist.

DEFINITION. A *Morse-Smale function* is a Morse function, $f : \mathbb{M} \rightarrow \mathbb{R}$, whose stable and unstable manifolds intersect transversally.

Roughly, this requires that the stable and unstable manifolds cross when they intersect. More formally, let $\sigma : \mathbb{R}^q \rightarrow \mathbb{M}$ and $v : \mathbb{R}^p \rightarrow \mathbb{M}$ be two smooth maps. Letting $z \in \mathbb{M}$ be a point in their common image, we say that σ and v intersect *transversally* at z if the derived images of the tangent spaces at preimages $x \in \sigma^{-1}(z)$ and $y \in v^{-1}(z)$ span the entire tangent space of \mathbb{M} at z :

$$D\sigma_x(T\mathbb{R}_x^q) + Dv_y(T\mathbb{R}_y^p) = T\mathbb{M}_z.$$

We say that σ and v are *transversal* to each other if they intersect transversally at every point z in their common image.

Complexes. Assuming transversality, the intersection of a stable q -manifold and an unstable p -manifold has dimension $q+p-d$. Furthermore, the boundary of every stable manifold is a union of stable manifolds of lower dimension. The set of stable manifolds thus forms a complex which we construct one dimension at a time.

0-skeleton: add all minima as stable 0-manifolds to initialize the complex;

1-skeleton: add all stable 1-manifolds, each an open interval glued at its endpoints to two points in the 0-skeleton;

2-skeleton: add all stable 2-manifolds, each an open disk glued along its boundary circle to a cycle in the 1-skeleton;

etc. It is possible that the two minima are the same so that the interval whose ends are both glued to it forms a loop. Similarly, the cycle in the 1-skeleton can be degenerate, such as pinched or even just a single point. Similar situations are possible for higher-dimensional stable manifolds. An example is the height function of the d -sphere. It has a single minimum, a single maximum, and no other critical points. The minimum has index 0 and forms a vertex in the complex. The maximum has index d and defines a stable d -manifold. It wraps around the sphere, and its boundary is glued to a single point, the minimum, as illustrated for $d = 2$ in Figure VI.6.

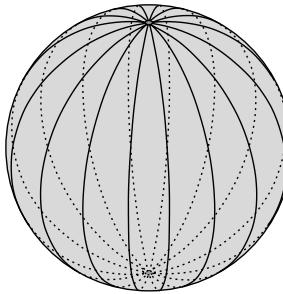


Figure VI.6: All integral lines of the height function of \mathbb{S}^2 originate at the minimum and end at the maximum. We therefore have two stable manifolds, a vertex for the minimum and an open disk for the maximum.

Morse inequalities. If we take the alternating sum of the numbers of stable manifolds in the above example, we get $1 + (-1)^d$, which is the Euler characteristic of the d -sphere. This is not a coincidence. More generally, the alternating sum of the numbers of stable manifolds gives the Euler characteristic, and this equation is one of the strong Morse inequalities. We state both the weak and the strong Morse inequalities, writing c_q for the number of critical points of index q .

MORSE INEQUALITIES. Let \mathbb{M} be a manifold of dimension d and let $f : \mathbb{M} \rightarrow \mathbb{R}$ be a Morse function. Then

- (i) **WEAK:** $c_q \geq \beta_q(\mathbb{M})$ for all q ;
- (ii) **STRONG:** $\sum_{q=0}^j (-1)^{j-q} c_q \geq \sum_{q=0}^j (-1)^{j-q} \beta_q(\mathbb{M})$ for all j .

As mentioned above, the strong Morse inequality for $j = d$ is an equality. We can recover the weak inequalities from the strong ones. Indeed

$$\begin{aligned} \sum_{q=0}^j (-1)^{j-q} c_q &\geq \beta_j(\mathbb{M}) - \sum_{q=0}^{j-1} (-1)^{j-q-1} \beta_q(\mathbb{M}) \\ &\geq \beta_j(\mathbb{M}) - \sum_{q=0}^{j-1} (-1)^{j-q-1} c_q. \end{aligned}$$

Removing the common terms on both sides leaves $c_j \geq \beta_j(\mathbb{M})$, the j -th weak inequality. We omit the proof of the strong inequalities and instead refer to the proof of their piecewise linear versions in the next section.

Floer homology. Assuming a Morse-Smale function, we can intersect the stable and unstable manifolds and get a refinement of the two complexes, which we refer to as the *Morse-Smale complex* of f . Its vertices are the critical points, and its cells are the components of the unions of integral lines with common origin and

common destination. It is quite possible that the stable manifold of a critical point intersects the unstable manifold of another critical point in more than one component. By definition of transversality, the index difference between the origin and the destination equals the dimension of the cell. In particular, the edges are isolated integral lines connecting index $q - 1$ with index q critical points.

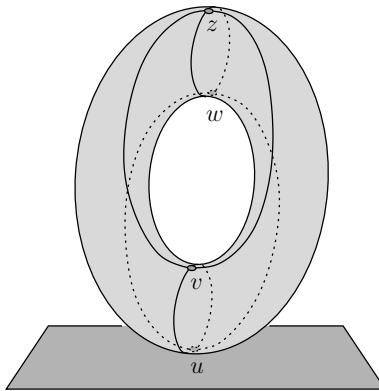


Figure VI.7: The Morse-Smale complex of the height function for the almost but not entirely upright torus.

To recover the homology of the manifold, we set up a chain complex. The q -chains are the formal sums of index q critical points. The boundary of an index q critical point, u , is the sum of index $q - 1$ critical points connected to u by an edge in the Morse-Smale complex. If there are multiple edges, we add the index $q - 1$ point multiple times. We illustrate this construction with the example depicted in Figure VI.7. We have a slightly tilted torus whose height function is a Morse-Smale function. There are one minimum, two saddles, and one maximum. The non-trivial chain groups are therefore $C_0 \simeq G$, $C_1 \simeq G^2$, $C_2 \simeq G$, with $G = \mathbb{Z}_2$, as usual. In this example, each critical point appears twice in the boundary of every other critical point, or not at all. Hence, the boundary of each one of the four critical points is zero. It follows that the boundary groups are trivial and the cycle groups as well as the homology groups are isomorphic to the chain groups. The Betti numbers are therefore $\beta_0 = 1$, $\beta_1 = 2$, $\beta_2 = 1$, which is consistent with what we already know about the torus.

Bibliographic notes. The concepts of integral lines and stable as well as unstable manifolds rely on fundamental properties of solutions to ordinary differential equations, in particular the Existence and Uniqueness Theorems; see e.g. Arnold [10]. The extra requirement of transversality between stable and unstable manifolds that distinguishes Morse from Morse-Smale functions has been proven to be generic by Kupka [97] and Smale [133]. The chain complex whose groups are formal sums of critical points is sometimes referred to as the Morse-Smale-Witten complex and the resulting homology theory is referred to as Floer homology [69].

VI.3 Piecewise Linear Functions

We rarely find smooth functions in practical situations. Instead, we often find non-smooth functions that approximate smooth functions or a series of non-smooth functions that approach a smooth limit. In this section, we turn things around and use insights gained into the smooth case as a guide in our attempt to understand the piecewise linear case.

Lower star filtration. Let K be a simplicial complex with real values specified at all vertices. Using linear extension over the simplices, we obtain a *piecewise linear (PL) function* $f : |K| \rightarrow \mathbb{R}$. It is defined by $f(x) = \sum_i b_i(x)f(u_i)$, where the u_i are the vertices of K and the $b_i(x)$ are the barycentric coordinates of x ; see Section III.1. It is convenient to assume that f is *generic*, by which we mean that the vertices have distinct function values. We can then order the vertices by increasing function value as $f(u_1) < f(u_2) < \dots < f(u_n)$. For each $0 \leq i \leq n$, we let K_i be the full subcomplex defined by the first i vertices. In other words, a simplex $\sigma \in K$ belongs to K_i iff each vertex u_j of σ satisfies $j \leq i$. Recall that the star of a vertex u_i is the set of cofaces of u_i in K . The *lower star* is the subset of simplices for which u_i is the vertex with maximum function value:

$$\text{St}_- u_i = \{\sigma \in \text{St } u_i \mid x \in \sigma \Rightarrow f(x) \leq f(u_i)\}.$$

Similar to the star, the lower star is generally not a complex. Adding the missing faces to the set, we get the *closed lower star*, which is a subcomplex of K . By assumption of genericity, each simplex has a unique maximum vertex and thus belongs to a unique lower star. It follows that the lower stars partition K . Furthermore, K_i is the union of the first i lower stars. This motivates us to call the nested sequence of complexes $\emptyset = K_0 \subseteq K_1 \subseteq \dots \subseteq K_n = K$ the *lower star filtration* of f . It will be useful to notice that the K_i are representative of the continuous family of sublevel sets. Specifically, for $f(u_i) \leq a < f(u_{i+1})$ the sublevel set $|K|_a = f^{-1}(-\infty, a]$

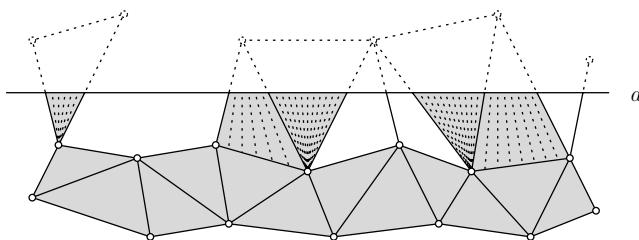


Figure VI.8: We retract $|K|_a$ to $|K_i|$ by shrinking the line segments decomposing the partial simplices from the top downward.

has the same homotopy type as K_i . To prove this, consider each simplex with at least one vertex in K_i and at least one vertex in $K - K_i$. Write this simplex as a union of line segments connecting points on the maximal face in K_i with points on the maximal face in $K - K_i$. In other words, express the simplex as the join of

these two faces; see Figure VI.8. The sublevel set contains only a fraction of each line segment, namely the portion from the lower endpoint x in $|K_i|$ to the upper endpoint y with $f(y) = a$. To get a deformation retraction, we let $(1-t)y + tx$ be the upper endpoint at time t . Going from time $t = 0$ to $t = 1$ deformation retracts $|K|_a$ onto $|K_i|$, so they have the same homotopy type.

PL critical points. We study the change from one complex to the next in the lower star filtration in more detail. Recall that the link of a vertex is the set of simplices in the closed star that do not belong to the star. Similarly, the *lower link* is the collection of simplices in the closed lower star that do not belong to the lower star. Equivalently, it is the collection of simplices in the link whose vertices have smaller function value than u_i :

$$\text{Lk}_- u_i = \{\sigma \in \text{Lk } u_i \mid x \in \sigma \Rightarrow f(x) < f(u_i)\}.$$

When we go from K_{i-1} to K_i , we attach the closed lower star of u_i , gluing it along the lower link to the complex K_{i-1} . Assume now that K triangulates a d -manifold. This restricts the possibilities dramatically since every vertex star is an open d -ball and every vertex link is a $(d-1)$ -sphere. A few examples of lower stars and lower links in a 2-manifold are shown in Figure VI.9. We classify the vertices using the

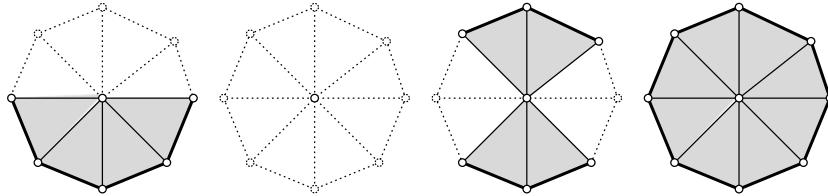


Figure VI.9: From left to right: the lower star and lower link of a regular vertex, a minimum, a saddle, and a maximum.

reduced Betti numbers of their lower links. Recall that $\tilde{\beta}_0$ is one less than β_0 , the number of components. The only exception to this rule is the empty lower link, for which we have $\tilde{\beta}_0 = \beta_0 = 0$ and $\tilde{\beta}_{-1} = 1$. Table VI.1 gives the reduced Betti numbers of the lower links in Figure VI.9. We call u_i a *PL regular vertex* if its lower

	$\tilde{\beta}_{-1}$	$\tilde{\beta}_0$	$\tilde{\beta}_1$
regular	0	0	0
minimum	1	0	0
saddle	0	1	0
maximum	0	0	1

Table VI.1: Classification of the vertices in a PL function on a 2-manifold.

link is non-empty but homologically trivial, and we call u_i a *simple PL critical vertex* of *index* q if its lower link has the reduced homology of the $(q-1)$ -sphere. In other words, the only non-zero reduced Betti number of a simple PL critical vertex of

index q is $\tilde{\beta}_{q-1} = 1$. We call a piecewise linear function $f : |K| \rightarrow \mathbb{R}$ on a manifold a *PL Morse function* if (i) each vertex is either PL regular or simple PL critical and (ii) the function values of the vertices are distinct.

Unfolding. In contrast to the smooth case, PL Morse functions are not dense among the class of all PL functions. Equivalently, a PL function on a manifold may require a substantial perturbation before it becomes PL Morse. As an example, consider the piecewise linear version of a monkey saddle displayed in Figure VI.10. It is therefore not reasonable to assume a PL Morse function as input, but we can

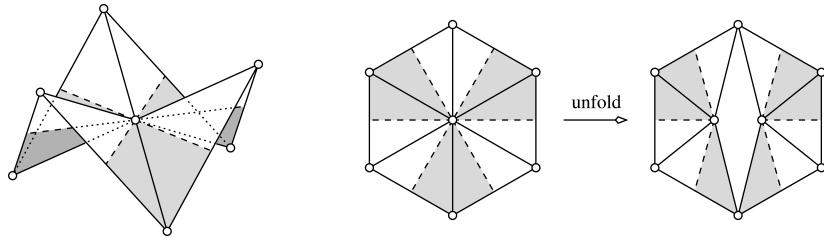


Figure VI.10: Left: a PL monkey saddle of a height function. The areas of points lower than the center vertex are shaded. Right: the unfolding of the monkey saddle into two simple saddles.

sometimes alter the triangulation locally to make it into a PL Morse function. In the 2-manifold case, a k -fold saddle is defined by $\tilde{\beta}_0 = k$. We can split it into k simple saddles by introducing $k - 1$ new vertices and assigning appropriate function values close to that of the original, k -fold saddle; see Figure VI.10 for the case $k = 2$. It is less clear how to unfold possibly complicated PL critical points for higher-dimensional manifolds.

Alternating sum of indices. Let K be a triangulation of a d -manifold and let $f : |K| \rightarrow \mathbb{R}$ be a PL Morse function. It is not difficult to prove that the alternating sum of the simple PL critical points gives the Euler characteristic:

$$\chi(K) = \sum_u (-1)^{\text{index}(u)}.$$

Since it is easy and instructive, we give an inductive proof of this equation. To go from K_{i-1} to K_i , we add the lower star of u_i . By the Euler-Poincaré Theorem, the Euler characteristic of the lower link, $A = \text{Lk}_- u_i$, is

$$\begin{aligned} \chi(A) &= \sum_{q \geq 1} (-1)^{q-1} \beta_{q-1}(A) \\ &= 1 + \sum_{q \geq 0} (-1)^{q-1} \tilde{\beta}_{q-1}(A). \end{aligned}$$

By definition, this is 1 if u_i is PL regular and $1 + (-1)^{\text{index}(u_i)-1}$ if u_i is PL critical. Each j -simplex in the lower star corresponds to a $(j-1)$ -simplex in the lower link,

except for the vertex u_i itself. Adding the lower star to the complex thus increases the Euler characteristic by $1 - \chi(A)$, which is zero for a PL regular point and $(-1)^{\text{index}(u_i)}$ for a simple PL critical point. The claimed equation follows.

Mayer-Vietoris sequences. We prepare the proof of the complete set of Morse inequalities for PL Morse functions by recalling the Mayer-Vietoris sequence of a covering of a simplicial complex by two subcomplexes. Let $K = K' \cup K''$ be the covering and note that the intersection of the two subcomplexes, $A = K' \cap K''$, is also a subcomplex of K . As discussed in Section IV.4, the reduced version of the corresponding Mayer-Vietoris sequence is

$$\dots \rightarrow \tilde{H}_{p+1}(K) \xrightarrow{\varphi_p} \tilde{H}_p(A) \xrightarrow{\psi_p} \tilde{H}_p(K') \oplus \tilde{H}_p(K'') \rightarrow \tilde{H}_p(K) \rightarrow \tilde{H}_{p-1}(A) \rightarrow \dots$$

It is exact, which means that the image of every homomorphism is equal to the kernel of the next homomorphism in the sequence. We are interested in the reduced

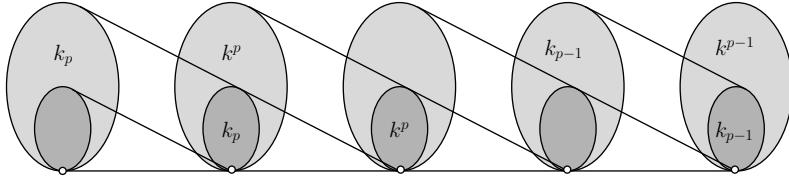


Figure VI.11: A portion of the Mayer-Vietoris sequence. By exactness, the rank of the kernel of every map complements the rank of the cokernel of the preceding map.

p -th homology group of A . Let k_p be the rank of $\ker \psi_p$ and let k^p be the rank of $\text{cok } \varphi_p = \tilde{H}_p(A)/\text{im } \varphi_p$. By exactness at $\tilde{H}_p(A)$, we have $\beta_p(A) = k_p + k^p$. As illustrated in Figure VI.11, exactness also implies that the rank of $\text{im } \psi_p$ is k^p and the rank of $\tilde{H}_{p+1}(K)/\ker \varphi_p$ is k_p .

We note that $\ker \psi_p$ and $\text{cok } \varphi_p$ distinguish two kinds of cycles in A . A cycle in the kernel bounds both in K' and in K'' , and these two $(p+1)$ -chains fit together to make a cycle of dimension $p+1$ in K . In contrast, a cycle in the cokernel is not in the image of the connecting homomorphism and thus represents a non-trivial homology class in K' or in K'' or in both.

PL Morse inequalities. We are now ready to state and prove the PL versions of the weak and strong Morse inequalities.

PL MORSE INEQUALITIES. Let K be a triangulation of a manifold of dimension d and let $f : |K| \rightarrow \mathbb{R}$ be a PL Morse function. Writing c_q for the number of index q PL critical points of f , we have

- (i) **WEAK:** $c_q \geq \beta_q(K)$ for all q ;
- (ii) **STRONG:** $\sum_{q=0}^j (-1)^{j-q} c_q \geq \sum_{q=0}^j (-1)^{j-q} \beta_q(K)$ for all j .

PROOF. We prove the inequalities inductively, for each K_i . They hold initially, when K_0 is empty. For the inductive step, we note that K_i is the union of K_{i-1} and the closed lower star of u_i . To study the situation, we use the Mayer-Vietoris sequence obtained by setting $K = K_i$, $K' = K_{i-1}$, $K'' = \text{St}_- u_i \cup \text{Lk}_- u_i$, and $A = \text{Lk}_- u_i$. Since K'' is the cone over a complex, it is homologically trivial. Referring to Figure VI.11, we let $\varphi_p : \tilde{H}_{p+1}(K) \rightarrow \tilde{H}_p(A)$ be the connecting homomorphism and $\psi_p : \tilde{H}_p(A) \rightarrow \tilde{H}_p(K') \oplus \tilde{H}_p(K'')$ be induced by inclusion. Furthermore, $k_p = \text{rank } \ker \psi_p$ and $k^p = \text{rank } \text{cok } \varphi_p$, as before. Since K'' is homologically trivial, the rank of $\tilde{H}_p(K)$ is the rank of $H_p(K')$ minus the rank of the image of ψ_p plus the rank of the kernel of ψ_{p-1} . Translating this back to the lower star filtration, we have

$$\text{rank } \tilde{H}_p(K_i) = \text{rank } \tilde{H}_p(K_{i-1}) - k^p + k_{p-1}.$$

By exactness of the sequence, $k_{p-1} + k^{p-1}$ is the rank of the reduced $(p-1)$ -st Betti number of A . This number is 1 if u_i is a simple PL critical point of index p and 0 otherwise. Specifically, if u_i is PL regular, then $k_{p-1} = k^{p-1} = 0$ for all p and the ranks of the homology groups do not change. Similarly, none of the counters of critical points change, so all Morse inequalities remain valid. If $\text{index}(u_i) = p$ and $k_{p-1} = 1$, then both c_p and $\tilde{\beta}_p$ go up by one, which maintains the validity of all Morse inequalities. On the other hand, if $\text{index}(u_i) = p$ and $k^{p-1} = 1$, then c_p goes up and $\tilde{\beta}_{p-1}$ goes down. Since the two have opposite signs, this maintains the validity of all Morse inequalities that contain both. The only strong Morse inequality that contains one but not both terms is the one for $j = p-1$. It contains the relevant Betti number with a plus sign, so this inequality is also preserved. \square

We note that the strong Morse inequality for $j = d$ is actually an equality, namely the one we have proved above, before recalling the Mayer-Vietoris sequence. It contains both changing terms, in all cases, so there is never a chance that the two sides become different. We also note that the proof of the Morse inequalities in the smooth case is the same. Indeed, passing a non-degenerate critical point has the same effect as adding the lower star of a simple PL critical vertex of the same index.

Bibliographic notes. Piecewise linear functions on polyhedral manifolds were studied by Banchoff [12]. He defines the index of a vertex as the Euler characteristic of its lower link. This is coarser than our definition but leads to similar results, in particular, a short and elementary proof that the Euler characteristic is equal to the alternating sum of critical points. However, it does not lend itself to a natural generalization of the other Morse inequalities to non-Morse PL functions. Our classification of PL critical points in terms of reduced Betti numbers can be found in [58], where it is used to compute the PL analog of the Morse-Smale complex for 2-manifolds. There are industrial applications of these ideas to surface design and segmentation based on curvature approximating and other shape-sensitive functions in \mathbb{R}^3 [55].

VI.4 Reeb Graphs

The structure of a continuous function can sometimes be made explicit by visualizing the evolution of the components of the level set. This leads to the concept of the Reeb graph of the function. It has applications in medical imaging and other areas of science and engineering.

Iso-surface extraction. The practical motivation for studying Reeb graphs is the extraction of iso-surfaces for 3-dimensional density data. In topological language, the density data is a continuous function, $f : [0, 1]^3 \rightarrow \mathbb{R}$, and an iso-surface is a level set, $f^{-1}(a)$. If f is smooth and a is a regular value of f as well as the restrictions of f to the faces of the cube, then the level set is a 2-manifold, possibly with boundary. Similarly, if f is generic PL and a is not the value of a vertex, then the level set is a 2-manifold, again possibly with boundary. Figure VI.12 illustrates this fact for a PL function on the unit square. Assuming we enter a triangle at

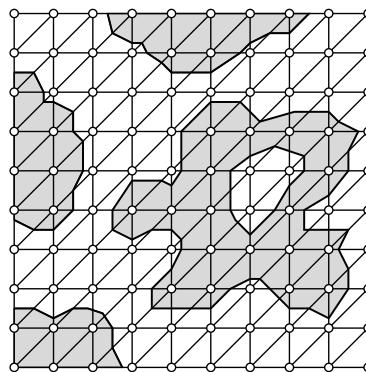


Figure VI.12: The level set of a generic PL function on a triangulation of the unit square. The superlevel set is white, and the sublevel set is shaded.

a boundary point x with $f(x) = a$, there is some other unique boundary point y with $f(y) = a$ where we exit the triangle. We draw the line segment from x to y as part of the level set and repeat the construction by entering the next triangle at y . There is never any choice as we trace the curve until we arrive at its other end. The procedure is similar for a PL function on the unit cube, except that we use a graph search algorithm to collect the triangular and quadrangular surface pieces we get by slicing the tetrahedra with planes. The most popular choices are Breadth-first Search and Depth-first Search, as described in Section II.2.

Given a first point on the level set, it is easy to trace out the component in which it is contained. However, to avoid missing any of the other components, we need to check the remainder of the triangulation. The desire to avoid this costly computation leads to the introduction of the contour tree. This is a data structure which can be queried for initial points on components of the level set without checking the entire triangulation. It is based on the concept of a Reeb graph.

Space of contours. Given a continuous map, $f : \mathbb{X} \rightarrow \mathbb{R}$, we note that the level sets form a partition of the topological space \mathbb{X} . We are interested in a possibly finer partition defined by calling two points $x, y \in \mathbb{X}$ *equivalent* if they belong to a common component of a level set of f . We refer to the equivalence classes as the *contours* of f . The *Reeb graph* of f is the set of contours, $R(f)$, together with the standard quotient topology. We recall that it is defined by taking all subsets whose preimages under $\psi : \mathbb{X} \rightarrow R(f)$ are open in \mathbb{X} , where $\psi(x)$ is of course the contour that contains x . Let $\pi : R(f) \rightarrow \mathbb{R}$ be the unique map whose composition with ψ is f . In other words, it is the map such that

$$\begin{array}{ccc} \mathbb{X} & \xrightarrow{f} & \mathbb{R} \\ & \psi \searrow & \nearrow \pi \\ & R(f) & \end{array}$$

commutes. We use it to explain how the Reeb graph speeds up the construction of a level set, $f^{-1}(a)$. Instead of going directly from \mathbb{R} to \mathbb{X} , we first compute the preimage of a under π , a set of points in the Reeb graph. The level set consists of a number of contours, one for each point r in $\pi^{-1}(a)$. In a medical imaging application, \mathbb{X} would be represented by a triangulation of the unit cube, and to go from a point r in $R(f)$ back to \mathbb{X} would be facilitated by a pointer to an edge in the triangulation that intersects the contour, $\psi^{-1}(r)$.

Besides using the Reeb graph as a data structure to accelerate the extraction of level sets, we may hope to learn something about the function or the topological space on which the function is defined. Even though the Reeb graph loses aspects of the original topological structure, there are some things it shows. First of all, $\psi : \mathbb{X} \rightarrow R(f)$ maps components to components. Furthermore, the Reeb graph reflects the 1-dimensional connectivity of the space in some cases. To describe this, we refer to a 1-cycle in $R(f)$ as a *loop* and write $\#\text{loops}$ for the size of the basis. The preimage of a loop in $R(f)$ is necessarily non-contractible in \mathbb{X} , and two different loops correspond to non-homologous 1-cycles. Expressing the two properties in terms of Betti numbers, we get

$$\begin{aligned} \beta_0(R(f)) &= \beta_0(\mathbb{X}), \\ \beta_1(R(f)) &\leq \beta_1(\mathbb{X}). \end{aligned}$$

Hence, if \mathbb{X} is contractible, then the Reeb graph is a tree, independent of the function f . In medical imaging, the space is a cube and thus contractible, which justifies the practice of calling $R(f)$ a contour tree.

Reeb graphs of Morse functions. More can be said if $\mathbb{X} = \mathbb{M}$ is a manifold of dimension $d \geq 2$ and $f : \mathbb{M} \rightarrow \mathbb{R}$ is a Morse function, as in Figure VI.13. Recall that each point $u \in R(f)$ is the image of a contour in \mathbb{M} . We call u a *node* of the Reeb graph if $\psi^{-1}(u)$ contains a critical point or, equivalently, if u is the image of a critical point under ψ . By definition of Morse function, the critical points

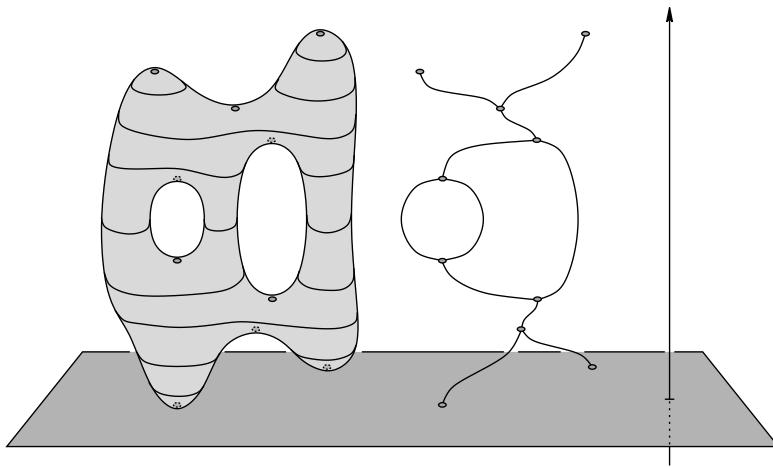


Figure VI.13: Level sets of the 2-manifold map to points on the real line and components of the level sets map to points of the Reeb graph.

have distinct function values, which implies a bijection between the critical points of f and the nodes of $R(f)$. The rest of the Reeb graph is partitioned into *arcs* connecting the nodes. A minimum starts a contour and therefore corresponds to a degree-1 node. An index-1 saddle that merges two contours into one corresponds to a degree-3 node. Symmetrically, a maximum corresponds to a degree-1 node and an index- $(d - 1)$ saddle that splits a contour into two corresponds to a degree-3 node. All other critical points correspond to nodes of degree two. Indeed, the only quadratic polynomials of the form $f(x) = -x_1^2 - \dots - x_q^2 + x_{q+1}^2 + \dots + x_d^2$ that have level sets with two components are the ones for $q = 1, d - 1$.

We note that the Reeb graph is a 1-dimensional topological space with points on arcs being individually meaningful objects. However, there is no preferred way to draw the graph in the plane or in space.

Loops in Reeb graphs. If \mathbb{M} is an orientable 2-manifold, then every saddle either merges two contours into one or splits a contour into two. Either way, the saddle corresponds to a degree-3 node in the Reeb graph. We use this fact to show that the number of loops depends only on \mathbb{M} and not on the function as long it is Morse. In the non-orientable case, we also have degree-2 nodes and therefore a number of loops that is no longer independent of the function.

LOOP LEMMA FOR 2-MANIFOLDS. The Reeb graph of a Morse function on a connected 2-manifold of genus g has g loops if the manifold is orientable and at most $\frac{g}{2}$ loops if it is non-orientable.

PROOF. Let c_q be the number of critical points of index q and n_i the number of nodes with degree i in the Reeb graph. We first consider the orientable case for

which the number of nodes is $n = n_1 + n_3$. We note that $n_1 = c_0 + c_2$ and $n_3 = c_1$. The number of arcs in the Reeb graph is $m = \frac{1}{2}(n_1 + 3n_3)$. The number of loops exceeds the surplus of arcs by one; that is,

$$\#\text{loops} = 1 + m - n = 1 - \frac{1}{2}(c_0 - c_1 + c_2).$$

By the last strong Morse inequality, the expression in parentheses is the Euler characteristic, which for orientable 2-manifolds is $\chi = 2 - 2g$. It follows that $\#\text{loops} = 1 - \frac{1}{2}(2 - 2g) = g$, as claimed. In the non-orientable case, the number of nodes is $n = n_1 + n_2 + n_3$, where $n_1 = c_0 + c_2$ and $n_2 + n_3 = c_1$. The number of arcs is $m = \frac{1}{2}(n_1 + 2n_2 + 3n_3)$. The number of loops is again one more than the surplus of arcs; that is,

$$\#\text{loops} = 1 + \frac{1}{2}(-n_1 + n_3) = 1 - \frac{1}{2}(c_0 - c_1 + c_2 + n_2).$$

Substituting the Euler characteristic for the alternating sum of critical points, we get $\#\text{loops} = 1 - \frac{1}{2}(\chi + n_2)$. For a non-orientable 2-manifold, we have $\chi = 2 - g$ and therefore $\#\text{loops} = \frac{1}{2}(g - n_2)$. Since the number of degree-2 nodes is non-negative, this is at most half the genus, as claimed. \square

Coincidentally, the proof implies that the number of degree-2 nodes has the same parity as the genus. Subject to this constraint, it can be anywhere between zero and g , which implies that the upper bound is tight and any integer number of loops between zero and half the genus can be achieved.

Constructing a Reeb graph. We finally consider the algorithmic problem of constructing the Reeb graph of a function on a 2-manifold. We assume the manifold is triangulated and the function, $f : M \rightarrow \mathbb{R}$, is PL Morse. The algorithm sweeps the manifold in the order of increasing function values. We thus begin by sorting the vertices such that $f(u_i) < f(u_{i+1})$ for $1 \leq i < n$. Consider a corresponding sequence of interleaved values, $s_1 < f(u_1) < s_2 < \dots < s_n < f(u_n) < s_{n+1}$. Since s_i is not the value of any vertex, its preimage is a 1-manifold, consisting of finitely many contours. Each contour is represented by a cyclic list of triangles in the triangulation. Every triangle contributes a line segment, and any two contiguous triangles meet in an edge that contributes a shared endpoint of two line segments to the contour. The representation is the same for all values strictly between $f(u_{i-1})$ and $f(u_i)$. Adjustments need to be made when we move into the next open interval, between $f(u_i)$ and $f(u_{i+1})$.

CASE 1: u_i is a minimum. Add a degree-1 node to the Reeb graph. It starts a new arc associated with a cyclic list initialized to the triangles in the star of u_i .

CASE 2: u_i is a regular vertex. Then two or more triangles in its star form a contiguous sequence in one of the cyclic lists. Except for the first and the last, all these triangles belong to the lower star. We remove the lower star triangles and replace them by the symmetrically defined upper star triangles of u_i .

CASE 3: u_i is a saddle. Then the triangles in its star form two contiguous sequences in the representation of the current level set. They may be part of the same cyclic list or of two different lists. Similarly to Case 2, we keep the first and last triangle of each sequence and replace the lower star triangles in between by the corresponding upper star triangles of u_i . Either list can be empty. We

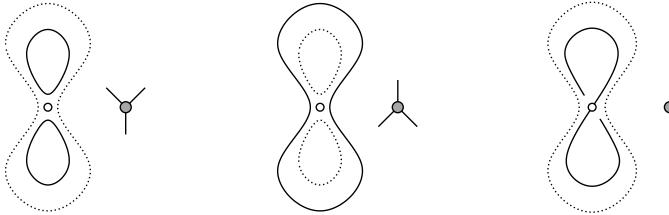


Figure VI.14: From left to right: merging two cyclic lists into one, splitting one list into two, reconnecting one list. Correspondingly, we add a down-fork, an up-fork, a degree-2 node to the Reeb graph.

do this by cutting the lists and regluing them when we add the upper star triangles. The global effect of the operation depends on whether the cutting is done on one or two cyclic lists and which ends are glued together. There are three different cases, as illustrated in Figure VI.14. In each case, we add a new node to the Reeb graph and represent the modified lists by arcs that end and start at that node.

CASE 4: u_i is a maximum. Remove the cyclic list of triangles in its star and end the corresponding arc by adding a new degree-1 node to the Reeb graph.

To implement the algorithm, we need a data structure that supports the following operations:

- CUT a cyclic list open by removing the links between two adjacent triangles;
- DROP a triangle from the end of an open list;
- APPEND a new triangle to the end of an open list;
- GLUE two ends of the same or of two different open lists;
- FIND the cyclic list that contains a specified triangle.

The cutting and gluing can be done without knowing whether the ends belong to the same or to different cyclic lists. However, to update the Reeb graph, we need to know which case we are in and we use the FIND operation to find out. All five operations are supported in time logarithmic in the length of the list if we store it in a data structure commonly referred to as a balanced search tree. Letting m be the number of edges in the triangulation, we thus get an algorithm that constructs the Reeb graph in time proportional to $m \log_2 m$. This is a significant improvement over the more straightforward algorithm that constructs the Reeb graph in time proportional to m^2 . No such improvement is currently known for functions on manifolds of dimension three or higher.

Bibliographic notes. The most common method for extracting iso-surfaces from density data is the Marching Cube Algorithm due to Lorensen and Cline [103]. As the name suggests, it works with a cube complex rather than a triangulation. The portion of the iso-surface within a single cube can be complicated, and the implementation of the algorithm requires some care. The idea of speeding up the iso-surface extraction with a contour tree is more recent [144]. This tree is really the Reeb graph of a PL function on a cube, which has no loops. The concept of the Reeb graph of a smooth function is much older [126]. The analysis of the number of loops and the Reeb Graph Algorithm for triangulated 2-manifolds are relatively recent results [39]. From a practical point of view, the most demanding operations are CUT and GLUE, as they require the splitting and melding of search trees. Particularly easy implementations of these operations are provided by the splay tree, a type of balanced search trees [132]. For contractible domains, the construction of the Reeb graph can be improved to time proportional to $m\alpha(m)$, where α is the extremely slow growing inverse of the Ackermann function [27].

Exercises

The credit assignment reflects a subjective assessment of difficulty. A typical question can be answered using knowledge of the material combined with some thought and analysis.

1. **Hessian** (two credits). Compute the Hessian and, if defined, the index of the origin, which is critical for each function in the list below.
 - (i) $f(x_1, x_2) = x_1^2 + x_2^2$.
 - (ii) $f(x_1, x_2) = x_1 x_2$.
 - (iii) $f(x_1, x_2) = (x_1 + x_2)^2$.
 - (iv) $f(x_1, x_2, x_3) = x_1 x_2 x_3$.
 - (v) $f(x_1, x_2, x_3) = x_1 x_2 + x_1 x_3 + x_2 x_3$.
 - (vi) $f(x_1, x_2, x_3) = (x_1 + x_2 + x_3)^2$.
2. **Approximate Morse function** (two credits). Let M be a geometrically perfect torus in \mathbb{R}^3 ; that is, M is swept out by a circle rotating about a line that lies in the same plane but does not intersect the circle. Let $f : M \rightarrow \mathbb{R}$ measure height parallel to the symmetry axis and note that f is not Morse.
 - (i) Describe a Morse function $g : M \rightarrow \mathbb{R}$ that differs from f by an arbitrarily small amount, $\|f - g\|_\infty < \varepsilon$.
 - (ii) Draw the Reeb graphs of both functions.
3. **Morse-Smale complex** (two credits). Let M be the torus in Exercise 2 and let $f : M \rightarrow \mathbb{R}$ measure height along a direction that is almost but not quite parallel to the symmetry axis of the torus.

- (i) Draw the Morse-Smale complex of the height function.
 - (ii) Give the chain, cycle, and boundary groups defined by Floer homology.
4. **Quadrangles** (three credits). Let \mathbb{M} be a 2-manifold and $f : \mathbb{M} \rightarrow \mathbb{R}$ a Morse-Smale function.
- (i) Prove that each 2-dimensional cell of the Morse-Smale complex of f is a quadrangle. In other words, each 2-dimensional cell is an open disk whose boundary can be decomposed into four arcs each glued to an edge in the complex.
 - (ii) Draw a case in which one edge is repeated so that the disk is glued to only three edges but twice to one of the three.
5. **Distance from a point** (three credits). Let \mathbb{M} be the torus swept out by a unit circle rotating at unit distance from the x_3 -axis. More formally, \mathbb{M} consists of all solutions to $x_1^2 + x_2^2 = (2 \pm \sqrt{1 - x_3^2})^2$ in \mathbb{R}^3 . For a point $z \in \mathbb{R}^3$ consider the function $f_z : \mathbb{M} \rightarrow \mathbb{R}$ defined by $f_z(x) = \|x - z\|$.
- (i) Describe the set of points z for which f_z violates the first property of a Morse function.
 - (ii) Describe the set of points z for which f_z is not a Morse function.
6. **Morse inequalities** (two credits). Recall that the unstable manifolds of a Morse function $f : \mathbb{M} \rightarrow \mathbb{R}$ are the stable manifolds of $-f$. Furthermore, if \mathbb{M} is a d -manifold, then an index- p critical point of f is an index- $(d-p)$ critical point of $-f$.
- (i) Use this symmetry to formulate collections of inequalities symmetric to the weak and strong Morse inequalities of f .
 - (ii) Use these inequalities to prove that the Euler characteristic of \mathbb{M} vanishes if d is odd.
7. **Reeb graph** (one credit). Consider the upright torus at time $t = 0$ and imagine it falling down in slow motion until it rests on its side at time $t = 1$.
- (i) What is the corresponding 1-parameter family of Reeb graphs of the height functions?
 - (ii) At which position (moment in time) does the Reeb graph not have a loop?
8. **BCC lattice** (two credits). Instead of the cubic lattice, we may consider constructing iso-surfaces from the body centered cubic (BCC) lattice obtained by adding the centers of all integer unit cubes. More formally, this is the set of points \mathbb{Z}^3 union $\mathbb{Z}^3 + (\frac{1}{2}, \frac{1}{2}, \frac{1}{2})^T$.
- (i) Show that there is an (infinite) simplicial complex whose vertex set is the BCC lattice and whose tetrahedra are pairwise congruent, that is, one can be obtained from any other by a rigid transformation.
 - (ii) Give a geometric description of the tetrahedron in (i), complete with all face, dihedral, and solid angles.

Part C

Computational Persistent Topology

Chapter VII

Persistence

The central concept of this chapter is motivated by the practical need to cope with noise in data. This includes defining, recognizing, and possibly eliminating noise. These are lofty goals and the challenge can be overwhelming. Indeed, the distinction between noise and feature is not well-defined but lies instead in the eye of the beholder. In any particular case, the focus is on a range of scales and a desire to ignore everything that is smaller or larger. In other words, we make ourselves the measure of all things and by doing so derive a unit, a point of view, and an opinion. Motivated by this thought, we take an agnostic approach and offer a means to measure scale, a tool that can be used to make judgments based on quantitative information, if one so desires.

VII.1 Persistent Homology

Persistent homology can be used to measure the scale or resolution of a topological feature. There are two ingredients, one geometric, defining a function on a topological space, and the other algebraic, turning the function into measurements. The measurements make sense only if the function does.

The elder rule. We begin with a simplified scenario in which we develop our intuition. Let \mathbb{X} be a connected topological space and $f : \mathbb{X} \rightarrow \mathbb{R}$ a continuous function. The sublevel sets of f form a 1-parameter family of nested subspaces, $\mathbb{X}_a \subseteq \mathbb{X}_b$ whenever $a \leq b$. It is convenient to write about this family as if it were one sublevel set that evolves as the threshold increases. We visualize this evolution by drawing each component of \mathbb{X}_a as a point. The result is a 1-dimensional graph, $G(f)$, not unlike the Reeb graph discussed in the previous chapter. Thinking of f as a height function, we draw the graph from bottom to top. Since components never shrink, the arcs of the graph may merge, but they never split. In the end, for large enough threshold a , we have a single component. It follows that $G(f)$ is a tree, and we refer to it as the *merge tree* of the function; see Figure VII.1.

We decompose this tree into disjoint paths that increase monotonically with f . To obtain the paths, we draw them from bottom to top, simultaneously, while keeping their upper endpoints at the same height, a . Paths extend; however, when they merge, we end the one that started later. Thinking of the difference between two function values as age, we give precedence to the older path.

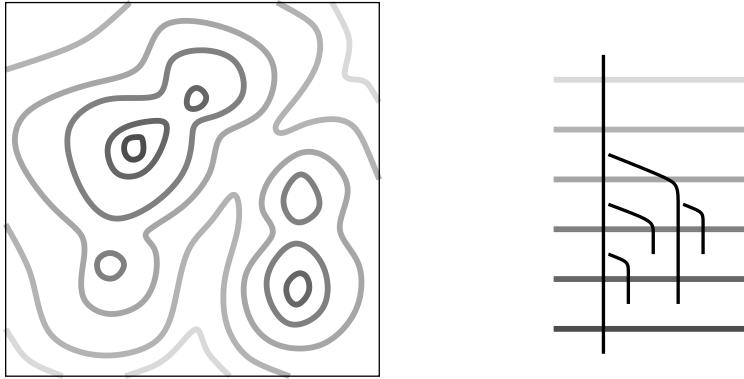


Figure VII.1: Left: a function on the unit square visualized by drawing six level sets with lighter shades of gray indicating larger values. Right: the path decomposition of the merge tree of the function.

ELDER RULE. At a juncture, the older of the two merging paths continues and the younger path ends.

Letting $a \leq b$ be two thresholds, we let $\beta(a, b)$ be the number of components in \mathbb{X}_b that have a non-empty intersection with \mathbb{X}_a . In terms of the merge tree, this is the number of subtrees with topmost points at value b that reach down to level a or below. Each such subtree has a unique path, its longest, that spans the entire interval between a and b . It follows that $\beta(a, b)$ is also the number of paths in the path decomposition of $G(f)$ that span $[a, b]$. We note that any path decomposition that is not generated using the Elder Rule does not have this property. In particular, if f is Morse, then the Elder Rule generates a unique path decomposition, which is the only one for which the number of paths spanning $[a, b]$ equals $\beta(a, b)$ for all values of $a \leq b$.

Filtrations. We obtain persistence by formulating the Elder Rule for the homology groups of all dimensions. Consider a simplicial complex, K , and a function $f : K \rightarrow \mathbb{R}$. We require that f be *monotonic*, by which we mean it is non-decreasing along increasing chains of faces, that is, $f(\sigma) \leq f(\tau)$ whenever σ is a face of τ . Monotonicity implies that the sublevel set, $K(a) = f^{-1}(-\infty, a]$, is a subcomplex of K for every $a \in \mathbb{R}$. Letting m be the number of simplices in K , we get $n+1 \leq m+1$ different subcomplexes, which we arrange as an increasing sequence:

$$\emptyset = K_0 \subseteq K_1 \subseteq \dots \subseteq K_n = K.$$

In other words, if $a_1 < a_2 < \dots < a_n$ are the function values of the simplices in K and $a_0 = -\infty$, then $K_i = K(a_i)$ for each i . We call this sequence of complexes the *filtration* of f and think of it as a construction by adding chunks of simplices at a time. We have seen examples before, namely the Čech and the alpha complexes in Chapter III and the lower star filtration of a piecewise linear function in Section VI.3. More than in the sequence of complexes, we are interested in the topological evolution, as expressed by the corresponding sequence of homology groups. For every $i \leq j$ we have an inclusion map from the underlying space of K_i to that of K_j and therefore an induced homomorphism, $f_p^{i,j} : H_p(K_i) \rightarrow H_p(K_j)$, for each dimension p . The filtration thus corresponds to a sequence of homology groups connected by homomorphisms,

$$0 = H_p(K_0) \rightarrow H_p(K_1) \rightarrow \dots \rightarrow H_p(K_n) = H_p(K),$$

again one for each dimension p . As we go from K_{i-1} to K_i , we might gain new homology classes and we might lose some when they become trivial or merge with each other. We collect the classes that are born at or before a given threshold and die after another threshold in groups.

DEFINITION. The p -th *persistent homology groups* are the images of the homomorphisms induced by inclusion, $H_p^{i,j} = \text{im } f_p^{i,j}$, for $0 \leq i \leq j \leq n$. The corresponding p -th *persistent Betti numbers* are the ranks of these groups, $\beta_p^{i,j} = \text{rank } H_p^{i,j}$.

Similarly, we define reduced persistent homology groups and reduced persistent Betti numbers. Note that $H_p^{i,i} = H_p(K_i)$. The persistent homology groups consist of the homology classes of K_i that are still alive at K_j or, more formally, $H_p^{i,j} = Z_p(K_j)/(B_p(K_j) \cap Z_p(K_i))$. We have such a group for each dimension p and each index pair $i \leq j$. We can be more concrete about the classes counted by the persistent homology groups. Letting γ be a class in $H_p(K_i)$, we say it is *born at K_i* if $\gamma \notin H_p^{i-1,i}$. Furthermore, if γ is born at K_i , then it *dies entering K_j* if it merges with an older class as we go from K_{j-1} to K_j , that is, $f_p^{i,j-1}(\gamma) \notin H_p^{i-1,j-1}$ but $f_p^{i,j}(\gamma) \in H_p^{i-1,j}$; see Figure VII.2. This is again the Elder Rule. If γ is born at K_i and dies entering K_j , then we call the difference in function value the *persistence*,

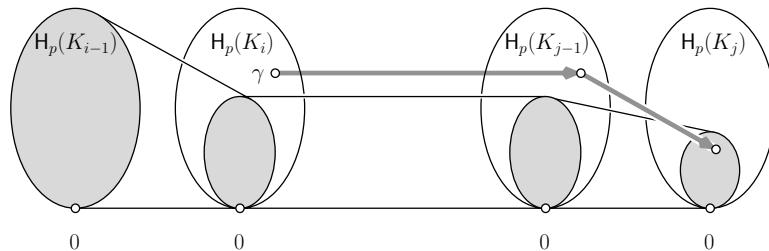


Figure VII.2: The class γ is born at K_i since it does not lie in the (shaded) image of $H_p(K_{i-1})$. Furthermore, γ dies entering K_j since this is the first time its image merges into the image of $H_p(K_{i-1})$.

$\text{pers}(\gamma) = a_j - a_i$. Sometimes we prefer to ignore the actual function values and consider the difference in index, $j - i$, which we call the *index persistence* of the class. If γ is born at K_i but never dies, then we set its persistence as well as its index persistence to infinity. We note that births and deaths can also be defined for a sequence of vector spaces that are not necessarily homology groups. All we need is a finite sequence and homomorphisms from left to right, which, for vector spaces, are usually referred to as linear maps.

Persistence diagrams. We visualize the collection of persistent Betti numbers by drawing points in two dimensions. Some of these points may have coordinates equal to infinity, and some might be the same, so we really talk about a multiset of points in the extended real plane, $\bar{\mathbb{R}}^2 = (\mathbb{R} \cup \{\pm\infty\})^2$. Letting $\mu_p^{i,j}$ be the number of independent p -dimensional classes that are born at K_i and die entering K_j , we have

$$\mu_p^{i,j} = (\beta_p^{i,j-1} - \beta_p^{i,j}) - (\beta_p^{i-1,j-1} - \beta_p^{i-1,j}),$$

for all $i < j$ and all p . Indeed, the first difference on the right-hand side counts the classes that are born at or before K_i and die entering K_j , while the second difference counts the classes that are born at or before K_{i-1} and die entering K_j . Drawing each point (a_i, a_j) with multiplicity $\mu_p^{i,j}$, we get the p -th *persistence diagram* of the filtration, denoted as $\text{Dgm}_p(f)$. It represents a class by a point whose vertical distance to the diagonal is the persistence. Since the multiplicities are defined only for $i < j$, all points lie above the diagonal. For technical reasons which will become clear in the next chapter, we add the points on the diagonal to the diagram, each with infinite multiplicity. Examples of persistence diagrams can be seen in Figure VII.5. It is easy to read off the persistent Betti numbers. Specifically, $\beta_p^{k,l}$ is the number of points in the upper left quadrant with corner point (a_k, a_l) . A class that is born at K_i and dies entering K_j is counted iff $a_i \leq a_k$ and $a_j > a_l$. The quadrant is therefore closed along its vertical right side and open along its horizontal lower side.

FUNDAMENTAL LEMMA OF PERSISTENT HOMOLOGY. Let $\emptyset = K_0 \subseteq K_1 \subseteq \dots \subseteq K_n = K$ be a filtration. For every pair of indices $0 \leq k \leq l \leq n$ and every dimension p , the p -th persistent Betti number is $\beta_p^{k,l} = \sum_{i \leq k} \sum_{j > l} \mu_p^{i,j}$.

This is an important property. It says the diagram encodes all information about persistent homology groups.

Matrix reduction. Besides having a compact description in terms of diagrams, persistence can also be computed efficiently. The particular algorithm we use is a version of matrix reduction. Perhaps surprisingly, we can get all the information with a single reduction. To describe this, we use a *compatible ordering* of the simplices, that is, a sequence $\sigma_1, \sigma_2, \dots, \sigma_m$ such that $i < j$ if $f(\sigma_i) < f(\sigma_j)$ or if σ_i is a face of σ_j . Such an ordering exists because f is monotonic. Note that every initial subsequence of simplices forms a subcomplex of K . We use this sequence

when we set up the m -by- m boundary matrix, ∂ , which stores the simplices of all dimensions in one place; that is,

$$\partial[i, j] = \begin{cases} 1 & \text{if } \sigma_i \text{ is a codimension-1 face of } \sigma_j; \\ 0 & \text{otherwise.} \end{cases}$$

In words, the rows and columns are ordered like the simplices in the total ordering and the boundary of a simplex is recorded in its column. The algorithm uses column operations to reduce ∂ to another 0-1 matrix R . Let $low(j)$ be the row index of the lowest 1 in column j . If the entire column is zero, then $low(j)$ is undefined. We call R reduced if $low(j) \neq low(j_0)$ whenever j and j_0 , with $j \neq j_0$, specify two non-zero columns. The algorithm reduces ∂ by adding columns from left to right.

```

 $R = \partial;$ 
 $\text{for } j = 1 \text{ to } m \text{ do}$ 
     $\text{while there exists } j_0 < j \text{ with } low(j_0) = low(j) \text{ do}$ 
         $\text{add column } j_0 \text{ to column } j$ 
     $\text{endwhile}$ 
 $\text{endfor.}$ 
```

The running time is at most cubic in the number of simplices. In matrix notation, the algorithm computes the reduced matrix as $R = \partial \cdot V$; see Figure VII.3. Since each simplex is preceded by its proper faces, ∂ is upper triangular. The j -th column of V encodes the columns in ∂ that add up to give the j -th column in R . Since we only add from left to right, V is also upper triangular and so is R .

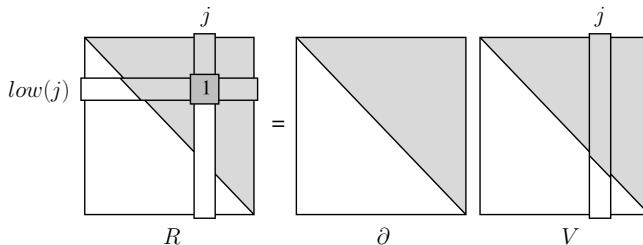


Figure VII.3: Reducing ∂ expressed as matrix multiplication. White areas are necessarily zero while entries in shaded areas can be either 0 or 1.

To get the ranks of the homology groups of K , we notice that the number of zero columns of R that correspond to p -simplices is the rank of Z_p . Similarly, the number of non-zero columns gives the rank of B_p . The difference is the p -th Betti number.

Pairing. However, there is significantly more information that we can harvest. To see this, we need to understand how the lowest 1s relate to the persistent homology groups. We begin by showing that they are unique, and this in spite of the fact that the reduced matrix, R , is not. Indeed, R is characterized by being reduced and

is obtained by left-to-right column operations. But we may or may not continue the operations once we have reached a reduced matrix. To see that the lowest 1s are unique, we consider the lower left submatrix R_i^j of R whose corner element is $R[i, j]$. In other words, R_i^j is obtained from R by removing the first $i - 1$ rows and the last $n - j$ columns. Since left-to-right column operations preserve the rank of every such submatrix, the rank of R_i^j is the same as that of the corresponding submatrix of ∂ , the one similarly obtained by removing the first $i - 1$ rows and the last $n - j$ columns. We consider the expression

$$r_R(i, j) = \text{rank } R_i^j - \text{rank } R_{i+1}^j + \text{rank } R_{i+1}^{j-1} - \text{rank } R_i^{j-1}$$

and note that $r_R(i, j) = r_\partial(i, j)$ for all i and j , where $r_\partial(i, j)$ has an analogous definition except when we take ranks of submatrices of ∂ . To evaluate this expression, we observe that the linear combination of any collection of non-zero columns in R_i^j is again non-zero. It follows that the rank of R_i^j is equal to its number of non-zero columns. Now, if $R[i, j]$ is a lowest 1, then R_i^j has one more non-zero column than the other three submatrices, which implies $r_R(i, j) = 1$. If $R[i, j]$ is not a lowest 1, then we consider two subcases. If none of the columns from 1 to $j - 1$ has its lowest 1 in row i , then R_i^j and R_{i+1}^j have the same number of non-zero columns and so do R_i^{j-1} and R_{i+1}^{j-1} . Second, if one of these columns has its lowest 1 in row i , then R_i^j has one more non-zero column than R_{i+1}^j and R_i^{j-1} has one more non-zero column than R_{i+1}^{j-1} . In either case, $r_R(i, j) = 0$. Since the ranks of the lower left submatrices of R are the same as those of ∂ , we have a characterization of the lowest 1s that does not depend on the reduction process.

PAIRING LEMMA. We have $i = \text{low}(j)$ iff $r_\partial(i, j) = 1$. In particular, the pairing between rows and columns defined by the lowest 1s in the reduced matrix does not depend on R .

Now that we know for sure that the lowest 1s are not an artifact of the particular strategy used for reduction, we ask what exactly they mean. Note that column j reaches its final form at the end of the j -th iteration of the outer loop. At this moment, we have the reduced matrix for the complex consisting of the first j simplices in the total ordering. We distinguish the case in which column j ends up zero from the other in which it has a lowest 1.

CASE 1: column j of R is zero. Consistent with the terminology introduced in Section V.4, we call σ_j positive since its addition creates a new cycle and thus gives birth to a new homology class.

CASE 2: column j of R is non-zero. It stores the boundary of the chain accumulated in column j of matrix V and is thus a cycle. Again consistent with the terminology in Section V.4, we call σ_j negative because its addition gives death to a homology class.

The class that dies in Case 2 is represented by column j . We still need to verify that it is born at the time the simplex of its lowest 1, σ_i with $i = \text{low}(j)$, is added. But

this is clear because the cycle in column j of R just died and all other cycles that die with it have 1s below row i ; otherwise, we could further reduce the matrix and obtain $\text{low}(j) < i$, which contradicts the algorithm. It follows that the lowest 1s indeed correspond to the points in the persistence diagrams. More precisely, (a_i, a_j) is a finite point in $\text{Dgm}_p(f)$ iff $i = \text{low}(j)$ and σ_i is a simplex of dimension p . In this case, σ_j is a simplex of dimension $p+1$. We have (a_i, ∞) in $\text{Dgm}_p(f)$ iff column i is zero but row i does not contain a lowest 1. In other words, σ_i is positive, but it does not get paired with a negative simplex.

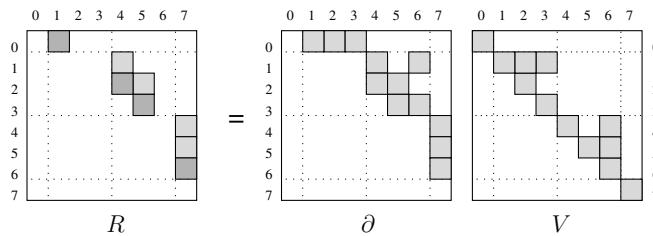


Figure VII.4: Reducing the boundary matrix of the complex consisting of a triangle and its faces. The shaded squares mark 1s in the matrices. The dark shaded squares mark lowest 1s in the reduced matrix.

An example. We illustrate the definitions with a small example. Let K consist of a triangle and its faces. To get a filtration, we first add the vertices, then the edges, and finally the triangle, numbering them in this order from 1 to 7. To make the exercise more interesting, we add the non-zero element of the (-1) -st reduced chain group as a dummy simplex of index 0 to compute reduced rather than ordinary homology. We recall that the augmentation map defines the boundary of each vertex as this dummy simplex. The resulting boundary matrix is shown as part of the matrix equation in Figure VII.4. We reduce it as described and get four non-

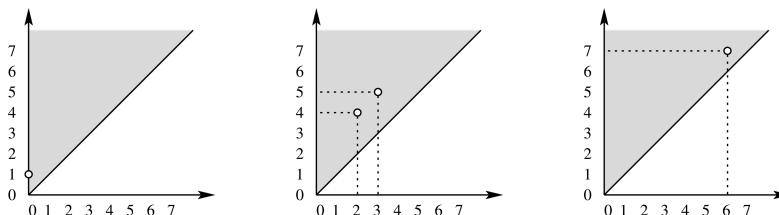


Figure VII.5: From left to right: the minus first, the zeroth, and the first persistence diagrams of the filtration that constructs a complex by first adding the three vertices, then the three edges, and finally the triangle.

zero columns in R . The first lowest 1 in R is in row 0 and column 1 and corresponds to the (-1) -dimensional reduced homology class that dies when we add vertex 1. The second lowest 1 is in row 2 and column 4. In words, the vertex 2 gives birth to the 0-cycle that the edge 4 kills. Similarly, the vertex 3 gives birth to the 0-cycle

that the edge 5 kills. Adding the edge 6 does not kill anything, which we see in the matrix since column 6 is zero. It corresponds to a 1-cycle obtained by adding the prior columns 4, 5, and 6, as indicated in V . The edge 6 thus gives birth to a 1-cycle that is then killed by the triangle 7. Figure VII.5 shows the corresponding three persistence diagrams which are drawn assuming the function value of a simplex is the same as its index. This particular function is injective, so all points in the diagrams have multiplicity one.

Bibliographic notes. The concept of persistent homology has been introduced for components by Frosini and Landi [73] and for general homology groups by Robins [127] and independently by Edelsbrunner, Letscher, and Zomorodian [60]. The latter paper gives the first fast algorithm for persistence, the same as described in this section but with the sparse matrix implementation discussed in the next section. A generalization of the notion of persistence to coefficient groups that are fields can be found in [161]; see also the monograph based on Zomorodian's thesis [160]. A recent survey on persistent homology is [57].

VII.2 Efficient Implementations

For practical applications, the number of simplices can be large so that storing the entire boundary matrix becomes prohibitive. As an alternative, we present a sparse matrix implementation of the Persistence Algorithm and give bounds on its running time that are better than cubic in the input size for many cases.

Sparse matrix representation. As in the previous section, we assume a monotonic function on a simplicial complex, $f : K \rightarrow \mathbb{R}$, and a compatible ordering of the simplices, $\sigma_1, \sigma_2, \dots, \sigma_m$. We store the data using a linear array, $\partial[1..m]$, and a linked list of simplices per entry. The list in $\partial[j]$ corresponds to the j -th column of the boundary matrix, storing the codimension-1 faces of σ_j . By the end of the algorithm, the list in the j -th array entry corresponds to the column of the reduced matrix whose lowest 1 is in the j -th row. If there is no such column, then the list will be empty. To emphasize the transition, we change the name for the array from ∂ at the beginning to R at the end of the algorithm. All lists are sorted in the order of decreasing index so that the most recently added simplex is readily available at the top; see Figure VII.6. We see a general migration of the lists from right to left.

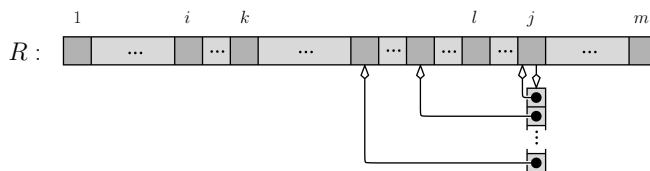


Figure VII.6: The sparse matrix representation of the reduced matrix with only one linked list shown.

To describe the algorithm that governs this migration, we write L for the linked list of the j -th array entry and $i = \text{TOP}(L)$ for the index of its top simplex. We call the i -th array entry *occupied* if it stores a non-empty list and *unoccupied* otherwise.

```

 $R = \partial;$ 
for  $j = 1$  to  $m$  do
     $L = \partial[j].cycle; R[j].cycle = \text{NULL};$ 
    while  $L \neq \text{NULL}$  and  $R[i]$  with  $i = \text{TOP}(L)$  is occupied do
         $L = L + R[i].cycle$ 
    endwhile;
    if  $L \neq \text{NULL}$  then  $R[i].cycle = L$  endif
endfor.
```

Adding two lists means merging them while deleting both copies of every duplicate simplex. Since we store the lists in consistent sorted order, each addition can be done in parallel scans. It is instructive to compare this sparse matrix version of the Persistence Algorithm with its standard matrix implementation.

Analysis. The main structure of the sparse matrix implementation is that of two nested loops, the outer and the inner loop. The addition of two lists is another loop in disguise, so the running time is at most cubic in the input size. To improve on this first estimate, we define a *collision* as an attempt to deposit the list L that fails because the entry is occupied. Each collision requires the merging of two lists, which takes time proportional to the sum of their lengths. The loop ends when L runs empty or when the non-empty list L is successfully deposited. The first case identifies σ_j as giving birth to a homology class. The second case identifies σ_j as giving death and the simplex, σ_i , where the deposit happens as triggering the corresponding birth. Each list $R[k].cycle$ contains σ_k as its topmost simplex. Similarly, σ_k is the topmost simplex in L when it collides with the list in $R[k]$. Using modulo 2 arithmetic, σ_k gets deleted, which implies that the topmost simplex in the merged list has index less than k . The inner loop thus proceeds monotonically from right to left. It follows that collisions for a simplex σ_j happen only at entries between i and j , where $i = 1$ if σ_j gives birth and i is the index of the corresponding birth if σ_j gives death. Note that in the latter case, $j - i$ is what we call the index persistence of σ_j . Consider now the inner loop for σ_j . A collision at entry k can happen only if σ_k gave birth to a class that died at σ_l before σ_j is reached. We have $i < k < l < j$, as in Figure VII.6. Similarly, the collisions during the inner loop for σ_l correspond to birth-death pairs nested within $[k, l]$. Inductively, this implies that the lists added at collisions contain only faces of simplices with index in $[i, j]$. Letting p be the dimension of σ_j , the number of such faces is at most $p + 1$ times the number of indices in the interval. The time to merge two lists is therefore at most proportional to this number. In summary, the running time of the inner loop for a p -simplex σ_j is at most $(p + 1)(j - i)^2$.

There are situations in which we know ahead of time which simplices give birth and which give death. For example, if the complex is geometrically realized in \mathbb{R}^3 , the Incremental Betti Number Algorithm described in Section V.4 gives such a

classification. With this information, we can then save the effort for the simplices that give birth so that the total running time of the algorithm becomes output-sensitive, and in particular bounded by the dimension times the sum of squares of the index persistences. Assuming constant dimension, this is at most proportional to m^3 , but for most practical data it is significantly smaller than that.

Zeroth diagram. The structure of the lists used to compute the 0-th persistence diagram is simpler than for dimensions beyond zero. This diagram depends solely on the vertices and edges of K and on their sequence in the compatible ordering. A vertex has no boundary and always gives birth to a component, so no choice there. An edge σ_j has two vertices as its boundary, $\partial\sigma_j = u + w$. Suppose u comes first, that is, $u = \sigma_i$, $w = \sigma_k$, and $i < k$. The first step of the algorithm is then its attempt to deposit the list L consisting of u and w in $R[k]$. If $L_k = R[k].cycle$ is empty, then the deposit is successful, σ_k, σ_j is a pair, and the inner loop ends. Otherwise, L_k is itself a list of two vertices, v and w in which v comes first. Adding the two lists gives $L + L_k$, which consists of u and v . Indeed, all non-empty lists have length two so that each addition takes only constant time. This implies that the total effort for dimension 0 is at most the sum of indices, for edges that give birth, and at most the sum of index persistences, for edges that give death. In any case, this is bounded from above by m^2 .

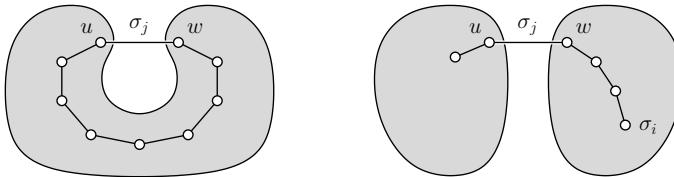


Figure VII.7: Adding the edge σ_j on the left gives birth to a 1-cycle while on the right it gives death to a component.

But we can do even better. Consider again the two cases for the edge with boundary $\partial\sigma_j = u + w$. It gives birth iff u and w belong to the same component of K_{j-1} , the complex right before we add σ_j ; see Figure VII.7 on the left. Starting with σ_j , the algorithm adds an edge to the growing path at each collision, and L keeps track of its boundary, the two endpoints. Eventually, the two ends meet, L becomes empty, and the path becomes a 1-cycle. The edge σ_j gives death iff u and w belong to two different components of K_{j-1} ; see Figure VII.7 on the right. The inner loop ends when one of the ends of the growing path reaches the first (oldest) vertex, σ_i , of one component. Since the inner loop works monotonically from right to left, this implies that the oldest vertex of the other component is even older. Following the Elder Rule, L gets deposited in $R[i]$ and σ_i, σ_j form a pair. Note that the outcome is predictable. All we need to know is whether or not u and w belong to different components in K_{j-1} , and if they do, which are the oldest vertices of these components. This is exactly the kind of information we can extract from the union-find data structure, as explained in Chapter I. Recall that this data structure stores each component as a tree of vertices. Given a vertex, we traverse the path up

to the root to determine the name of the component. Using the index of the oldest vertex as the name gives the information we need at negligible cost. In summary, we compute the 0-th persistence diagram in time at most proportional to $m\alpha(m)$, where α is the inverse of the Ackermann function which, for all practical purposes, is bounded from above by a constant.

Surfaces. We now consider a simplicial complex, K , that triangulates a 2-manifold. This case is of some practical importance and it allows for a fast implementation of the Persistence Algorithm. Let $f : |K| \rightarrow \mathbb{R}$ be obtained by piecewise linear interpolation of its values at the vertices, as explained in Section III.1. There is possibly non-trivial information in the 0-th and the 1-st persistence diagrams of f but not in any of the others. To compute these two diagrams fast, we need to answer two questions.

1. How can we turn the 1-parameter family of sublevel sets into a filtration that we can feed to our algorithm?
2. How can we improve the slower running time for the 1-st persistence diagram to roughly the time needed for the 0-th diagram?

We deal with the first question now and defer the second question to later. Assume for simplicity that the restriction of f to the vertices of K is injective. As defined in Chapter VI, the lower star filtration is then the sequence $\emptyset = K_0 \subseteq K_1 \subseteq \dots \subseteq K_n = K$, where K_i is the union of the lower stars of the first i vertices in the ordering by f . It is also the filtration generated by the monotonic function $g : K \rightarrow \mathbb{R}$ defined by mapping each simplex to $g(\sigma) = \max_{x \in \sigma} f(x)$. The diagrams of f are defined by the homology groups of the sublevel sets of f , $|K|_a = f^{-1}(-\infty, a]$, while those of g are defined by the homology groups of the sublevel sets of g , $K_a = g^{-1}(-\infty, a]$. By definition of lower star filtration, we have $|K_a| \subseteq |K|_a$, and the inclusion is a homotopy equivalence; see Figure VI.8 and the discussion around it. It follows that the vertical maps in the following diagram are isomorphisms:

$$\begin{array}{ccc} H_p(|K|_a) & \longrightarrow & H_p(|K|_b) \\ \uparrow & & \uparrow \\ H_p(K_a) & \longrightarrow & H_p(K_b), \end{array}$$

where p is any dimension and a, b , with $a \leq b$, are any two real numbers. The square commutes because all four maps are induced by inclusion. Indeed, these two conditions suffice for the diagrams defined by the two sequences to be the same.

PERSISTENCE EQUIVALENCE THEOREM. Consider two sequences of vector spaces connected by homomorphisms $\phi_i : U_i \rightarrow V_i$:

$$\begin{array}{ccccccc} V_0 & \rightarrow & V_1 & \rightarrow & \dots & \rightarrow & V_{n-1} \rightarrow V_n \\ \uparrow & & \uparrow & & & & \uparrow \\ U_0 & \rightarrow & U_1 & \rightarrow & \dots & \rightarrow & U_{n-1} \rightarrow U_n. \end{array}$$

If the ϕ_i are isomorphisms and all squares commute, then the persistence diagram defined by the U_i is the same as that defined by the V_i .

The proof is not difficult but it is tedious and is therefore omitted. As explained above, the 0-th persistence diagram of g can be computed in time at most proportional to $m\alpha(m)$. The equivalence with the 0-th persistence diagram of f thus implies that the latter can be computed in the same amount of time.

First diagram. Instead of computing the 1-st persistence diagram of f directly, we construct the 0-th persistence diagram of $-f$ and derive the diagram of f from it. We begin by describing the relation between $\text{Dgm}_1(f)$ and $\text{Dgm}_0(-f)$, omitting proofs since the relations are consequences of the more general theorems given in the next section.

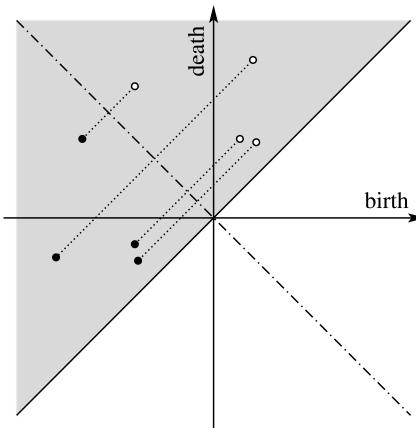


Figure VII.8: The white points of $\text{Dgm}_1(f)$ are reflections of the black points of $\text{Dgm}_0(-f)$ across the minor diagonal.

The 1-st persistence diagram of f consists of the diagonal, a finite portion of off-diagonal points (a, b) , and an infinite portion of off-diagonal points (c, ∞) . We construct the finite portion from the 0-th persistence diagram of $-f$. Specifically, the point (a, b) marks the birth of a 1-dimensional homology class at a and its death at b . Looking at $-f$ is like taking the complement and going backward. We thus have the birth of a 0-dimensional homology class at $-b$ and its death at $-a$. It follows that a point (a, b) belongs to $\text{Dgm}_1(f)$ iff the point $(-b, -a)$ belongs to $\text{Dgm}_0(-f)$. In other words, the finite portion of $\text{Dgm}_1(f)$ can be obtained by reflecting the finite portion of $\text{Dgm}_0(-f)$ across the minor diagonal, as illustrated in Figure VII.8. We get the points at infinity by partitioning the set of edges in the complex into three subsets: edges that give death in the lower star filtration of f , edges that give death in the lower star filtration of $-f$, and the rest. The first two contribute coordinates to the finite portions of the 0-th and the 1-st diagrams of f . For each edge in the third set, we have a point at infinity in the 1-st diagram, namely a class born when the edge is added and living on even when the complex K is complete. In summary, we have a three-pass algorithm for computing the persistence diagrams of a piecewise linear function f on a triangulated 2-manifold in time at most proportional to $m\alpha(m)$.

Bibliographic notes. The original paper on persistent homology by Edelsbrunner, Letscher, and Zomorodian [60] describes the sparse matrix version of the Persistence Algorithm explained in this section. Furthermore, the paper focuses on cases in which birth and death information is available using the Incremental Betti Number Algorithm by Delfinado and Edelsbrunner [45]. The standard matrix reduction version of the Persistence Algorithm came later historically and brought with it a more general appeal at the expense of increased computational resources. The Persistence Equivalence Theorem relating diagrams of different functions first appeared in [161].

VII.3 Extended Persistence

In this section, we discuss an extension of persistence that is motivated by an approach to fitting shapes to each other. The problem of fitting shapes arises when we solve a puzzle but also in the assembly of mechanical shapes, in the reconstruction of broken artifacts, and in protein docking.

Elevation on a surface. We give a brief sketch of the approach to fitting shapes and refer to Section IX.2 for a more detailed description. Let \mathbb{M} be a smoothly embedded 2-manifold in \mathbb{R}^3 . Given a direction $u \in \mathbb{S}^2$, the *height function* in this direction, $f : \mathbb{M} \rightarrow \mathbb{R}$, is defined by mapping each point x to $f(x) = \langle x, u \rangle$. We usually draw u vertically going up and think of the height as the signed distance from a horizontal base plane, as in Figure VII.9. Given a threshold $a \in \mathbb{R}$, we recall

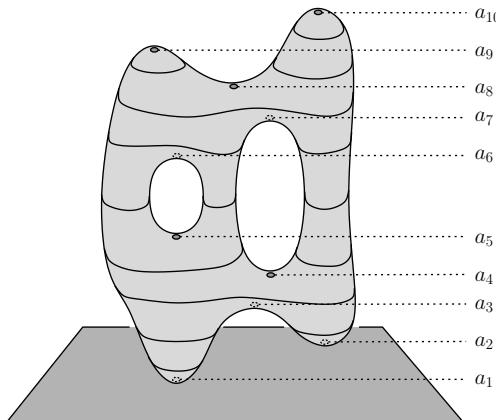


Figure VII.9: A smoothly embedded 2-manifold with level sets shown and critical points of the vertical height function marked.

that the sublevel set consists of all points with height a or less, $\mathbb{M}_a = f^{-1}(-\infty, a]$. As mentioned in the previous sections, the sublevel sets are nested and define persistence through the corresponding sequence of homology groups. For a generic

smooth surface, the homological critical values of a height function are the height values of isolated critical points. If, furthermore, the direction is generic, then there are only three different types: minima which start components, saddles which merge components or complete loops, and maxima which fill holes. Assuming the critical points have distinct heights, the points in the persistence diagrams of f correspond to pairs of critical points. The elevation at the points x and y of such a pair is set to $|f(x) - f(y)|$. Since x is critical for two opposite directions, we need to make sure that the pairing is the same in both directions, else we get contradictory assignments of elevation. We also need all critical points to be paired; otherwise, we get white areas in which elevation remains undefined. The latter is the reason for why we extend persistence and the former is a constraint we need to observe in this extension.

Extended filtration. Let $a_1 < a_2 < \dots < a_n$ be the homological critical values of the height function $f : M \rightarrow \mathbb{R}$. At interleaved values

$$b_0 < a_1 < b_1 < a_2 < \dots < a_n < b_n$$

we get sublevel sets $M_{b_i} = f^{-1}(-\infty, b_i]$ which are 2-manifolds with boundary. Symmetrically, we define *superlevel sets* $M^{b_i} = [b_i, \infty)$, which are complementary 2-manifolds with the same boundary. Finally, we use both to construct a sequence of homology groups going up and a sequence of relative homology groups coming back down:

$$\begin{aligned} 0 &= H_p(M_{b_0}) &\rightarrow \dots \rightarrow H_p(M_{b_n}) \\ &= H_p(M, M^{b_n}) &\rightarrow \dots \rightarrow H_p(M, M^{b_0}) &= 0 \end{aligned}$$

for each dimension p . The homomorphisms are induced by inclusion. We recall that for modulo 2 arithmetic, the homology groups are isomorphic to the cohomology groups. Furthermore, Lefschetz duality implies $H^p(M_b) \cong H_{d-p}(M, M^b)$. This shows that the construction is intrinsically symmetric although not necessarily within the same dimension. Since we go from the trivial group to the trivial group, everything that gets born eventually dies. As a consequence, all births will be paired with corresponding deaths, as desired.

Tracing what gets born and dies in the relative homology groups is a bit less intuitive than for the absolute homology groups going up. However, we can translate the events between the absolute homology of M^b and the relative homology of the pair (M, M^b) . Coming down, the threshold decreases, so the superlevel set grows. We call a homology class in the superlevel set *essential* if it lives all the way down to b_0 and *inessential* otherwise.

RULE 1: a dimension p homology class of M^b dies at the same time that a dimension $p + 1$ relative homology class of (M, M^b) dies.

RULE 2: an inessential dimension p homology class of M^b gets born at the same time that a dimension $p + 1$ relative homology class of (M, M^b) gets born.

RULE 3: an essential dimension p homology class of M^b gets born at the same time that a dimension p relative homology class of (M, M^b) dies.

We can prove these relationships by studying the kernels and cokernels of the maps from the homology groups of \mathbb{M}^b into those of \mathbb{M} . Leaving this to the interested reader, we develop our intuition by considering an example.

Example. Consider the height function of the genus-2 torus in Figure VII.9. Going up, a_1 and a_2 give birth to classes in H_0 , a_4, a_5, a_6, a_7, a_8 give birth to classes in H_1 , and a_{10} gives birth to a class in H_2 . All classes live until the end of the ascending pass, except for the dimension 0 class born at a_2 , which dies at a_3 , and the dimension 1 class born at a_8 , which dies at a_9 . These are the only two finite off-diagonal points in the ordinary persistence diagrams. Coming down, a_{10} kills the class in H_0 and a_9 gives birth to a class in H_1 that dies at a_8 . Furthermore, a_7, a_6, a_5, a_4 kill the classes in H_1 , a_3 gives birth to a class in H_2 that dies at a_2 , and finally a_1 kills the class in H_2 that was born going up at a_{10} . To summarize, the pairs of critical values defining the points in the diagrams are $(a_1, a_{10}), (a_2, a_3)$ in dimension 0, $(a_4, a_7), (a_5, a_6), (a_6, a_5), (a_7, a_4), (a_8, a_9), (a_9, a_8)$ in dimension 1, and $(a_{10}, a_1), (a_3, a_2)$ in dimension 2. We show the diagrams in Figure VII.10 using different symbols for classes born and dying going up, born going up and dying coming down, and born and dying coming down. They make up the *ordinary*, the

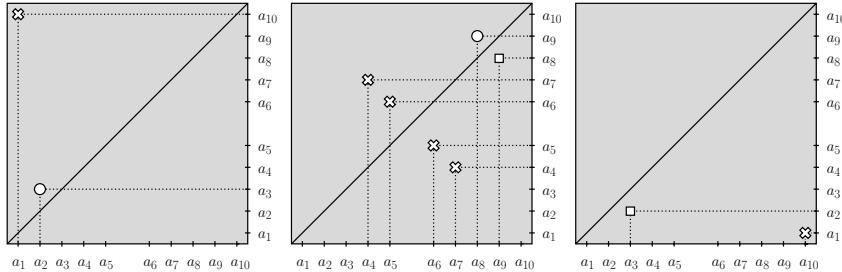


Figure VII.10: From left to right: the 0-th, 1-st, 2-nd persistence diagrams of the height function in Figure VII.9.

extended, and the *relative subdiagrams*, which we denote as Ord, Ext, and Rel, with the dimension in the index and the function in parentheses, as before. Note that the points of the ordinary subdiagrams lie above and those of the relative subdiagrams lie below the diagonal. The points of the extended subdiagrams can lie on either side.

Duality and symmetry. The symmetries we observe in Figure VII.10 are not coincidental. They arise as consequences of Lefschetz duality between absolute and relative homology groups of complementary dimensions, $H_p(\mathbb{M}_b) \simeq H_{d-p}(\mathbb{M}, \mathbb{M}^b)$. This translates into a duality result for persistence diagrams, which we state without proof. We use a superscript ‘ T ’ to indicate reflection across the main diagonal, mapping the point (a, b) to (b, a) .

PERSISTENCE DUALITY THEOREM. Let f be a function on a d -manifold without boundary. Then the persistence diagrams are reflections of each other as follows:

$$\begin{aligned}\text{Ord}_p(f) &= \text{Rel}_{d-p}^T(f), \\ \text{Ext}_p(f) &= \text{Ext}_{d-p}^T(f), \\ \text{Rel}_p(f) &= \text{Ord}_{d-p}^T(f).\end{aligned}$$

Equivalently, the full p -th persistence diagram is the reflection of the full $(d-p)$ -th persistence diagram, $\text{Dgm}_p(f) = \text{Dgm}_{d-p}^T(f)$. We have $d = 2$ for the example illustrated in Figures VII.9 and VII.10 and we indeed have diagrams that are reflections of each other as described. For $2p = d$, the extended subdiagram is the reflection of itself and is therefore symmetric across the main diagonal.

Recall that the definition of elevation requires that the pairing of critical points be the same for antipodal height functions. We can use duality to prove that they are indeed the same. More specifically, we have the following structural result, again expressed in terms of subdiagrams of the persistence diagrams. We use the superscript ‘ R ’ to indicate reflection across the minor diagonal, mapping the point (a, b) to $(-b, -a)$. Similarly, we use the superscript ‘ 0 ’ to indicate central reflection or rotation by 180 degrees, mapping the point (a, b) to $(-a, -b)$.

PERSISTENCE SYMMETRY THEOREM. Let f be a function on a d -manifold without boundary and let $-f$ be its negative. Then the persistence diagrams of the two functions are reflections of each other:

$$\begin{aligned}\text{Ord}_p(f) &= \text{Ord}_{d-p-1}^R(-f), \\ \text{Ext}_p(f) &= \text{Ext}_{d-p}^0(-f), \\ \text{Rel}_p(f) &= \text{Rel}_{d-p+1}^R(-f).\end{aligned}$$

In lieu of a proof, we just mention that each of the three equations can be obtained using the Persistence Duality Theorem together with the above three rules relating events in the parallel sequences of absolute and relative homology groups.

Lower and upper stars. To describe how we compute extended persistence, let K be a triangulation of a d -manifold \mathbb{M} . We assume the height function is defined at the vertices. We also assume that the height values are distinct, so we can index the vertices such that $f(v_1) < f(v_2) < \dots < f(v_n)$. Let $f : |K| \rightarrow \mathbb{R}$ be obtained by piecewise linear extension. Writing $a_i = f(v_i)$ and introducing interleaved values $b_0 < a_1 < b_1 < \dots < a_n < b_n$, we can define sublevel sets and superlevel sets as before. The set of points $x \in |K|$ with $f(x) \leq b_i$ is homeomorphic to \mathbb{M}_{b_i} and thus is a manifold with boundary. Similarly, the set of points with $f(x) \geq b_i$ is homeomorphic to \mathbb{M}^{b_i} and is a manifold with boundary. We can retract the partially used simplices and get homotopy equivalent subcomplexes of K . Specifically, let K_i be the full subcomplex defined by the first i vertices in the ordering and K^i the full subcomplex defined by the last $n - i$ vertices. The two subcomplexes of K

are disjoint although together they cover all n vertices. The only simplices not in either subcomplex are the ones that connect the first i with the last $n - i$ vertices. Recall that the lower star of a vertex v_i consists of all simplices that have v_i as their highest vertex. Symmetrically, we define the *upper star* to consist of all simplices that have v_i as their lowest vertex. More formally,

$$\begin{aligned} \text{St}_-v_i &= \{\sigma \in \text{St } v_i \mid x \in \sigma \Rightarrow f(x) \leq f(v_i)\}, \\ \text{St}^+v_i &= \{\sigma \in \text{St } v_i \mid x \in \sigma \Rightarrow f(x) \geq f(v_i)\}. \end{aligned}$$

Since every simplex has a unique highest and a unique lowest vertex, the lower stars partition K and so do the upper stars. With this notation, $K_0 = \emptyset$ and $K_i = K_{i-1} \cup \text{St}_-v_i$ for $1 \leq i \leq n$. Equivalently, K_i is the union of the first i lower stars. Symmetrically, $K^n = \emptyset$, $K^i = K^{i+1} \cup \text{St}^+v_{i+1}$, and K^i is the union of the last $n - i$ upper stars.

Computation. By the Persistence Equivalence Theorem in the previous section, the K_i have the same homotopy type as the sublevel sets, and the K^i have the same homotopy types as the superlevel sets of \mathbb{M} . We can therefore use them to compute persistence. Let A be the boundary matrix for the ascending pass, storing the simplices in blocks that correspond to the lower stars of v_1 to v_n , in this order. Within each block, we store the simplices in order of non-decreasing dimension and break remaining ties arbitrarily. All simplices in the same block are assigned the same value, namely the height of the vertex defining the lower star. If two simplices in the same block are paired, they define a point on the diagonal of the appropriate persistence diagram. In other words, the homology class dies as soon as it is born and therefore has zero persistence. Only pairs between blocks carry any significance.

Let B be the boundary matrix for the descending pass, storing the simplices in blocks that correspond to the upper stars of v_n to v_1 , in this order. Using A and B , we form a bigger matrix by adding the zero matrix at the lower left and the permutation matrix P that translates between A and B at the upper right, as in Figure VII.11. We can think of the result as the boundary matrix of a new complex, namely the cone over K . We pick a new, dummy vertex, v_0 , and for each i -simplex σ in K add the $(i + 1)$ -simplex $\sigma \cup \{v_0\}$. Adding the cone removes any non-trivial homology. This explains why reducing the big matrix works. As we move from left to right, we first construct K , forming pairs by reducing A . At the halfway point, the only unpaired simplices are the ones that gave birth to the essential homology classes. As we continue, we cone off K step by step, eventually removing all non-trivial homology. In the end, the ordinary, extended, and relative subdiagrams are given by the lowest 1s in the upper left, upper right, and lower right quadrants of the reduced matrix.

Indeed, we draw the diagram that corresponds to one of the three quadrants by marking each lowest one as a point, replacing indices by function values. For A , the birth values increase downward and the death values from left to right, so we need to turn the quadrant by 90° to get the ordinary subdiagram. Symmetrically, we turn the quadrant of B by -90° to get the relative subdiagram and we reflect the quadrant of P across the main diagonal to get the extended subdiagram. Since

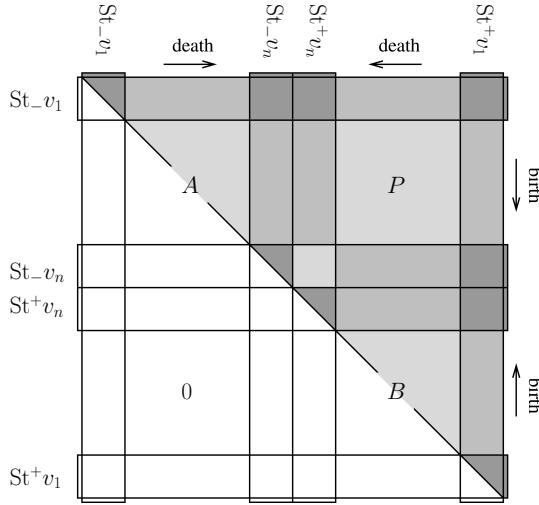


Figure VII.11: The block structure of the boundary matrix representing the construction of K going up and the subsequence destruction coming down.

the reduced versions of A and B are upper triangular, we indeed get the ordinary subdiagram above and the relative subdiagram below the diagonal.

Bibliographic notes. The extension of persistence described in this section is due to Cohen-Steiner, Edelsbrunner, and Harer [35]. It makes essential use of Poincaré and Lefschetz duality to obtain the desired symmetry properties for manifolds. The construction applies equally well to general topological spaces but without guarantee of duality and symmetry. The main motivation for the extension is the definition of the elevation function of a smoothly embedded surface in \mathbb{R}^3 ; see Section IX.2. This definition requires that all critical points be paired, which is not the case for ordinary persistence. The original paper on elevation contains an elementary description of extended persistence just for the case of surfaces [3].

VII.4 Spectral Sequences

Topologists will immediately recognize a connection between persistence and spectral sequences. We shed light on this relation by reviewing spectral sequences, first in terms of the matrix reduction algorithm and second in terms of groups and maps between them.

The matrix reduction view. As usual, we start with a filtration of a simplicial complex,

$$\emptyset = K_0 \subseteq K_1 \subseteq \dots \subseteq K_n = K,$$

letting $k_i = \text{card } K_i$ be the number of simplices in the i -th complex. Using a compatible total ordering of the simplices, we let ∂ be the boundary matrix which we write in block form. Specifically, ∂_i consists of the rows numbered $k_{i-1} + 1$ to k_i corresponding to the simplices in $K_i - K_{i-1}$, and ∂^j consists of the columns numbered $k_{j-1} + 1$ to k_j corresponding to the simplices in $K_j - K_{j-1}$. We write ∂_i^j for the intersection of the i -th block of rows and the j -th block of columns; that is, ∂_i^j records the codimension-1 faces of the simplices in $K_j - K_{j-1}$ that lie in $K_i - K_{i-1}$. Since the boundary matrix is upper triangular, we have $\partial_i^j = 0$ whenever $i > j$. We reduce the boundary matrix with left-to-right column additions, as before, but instead of sweeping the matrix from left to right, we sweep it diagonally. More precisely, we work in phases, and in Phase r , we reduce columns in ∂^j by adding columns in the blocks from ∂^{j-r+1} all the way to ∂^j itself. The Spectral Sequence Algorithm thus reduces the columns from the diagonal outward, as illustrated in Figure VII.12.

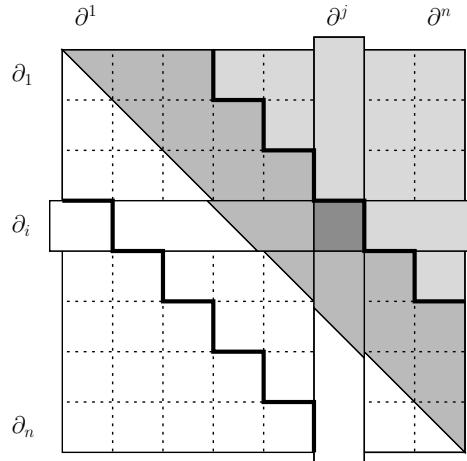


Figure VII.12: After three phases, the triple blocks along the diagonal are reduced. The highlighted blocks of rows and columns intersect in the block matrix ∂_i^j .

```

for r = 1 to n do
    for j = r to n do
        for i = k_{j-1} + 1 to k_j do
            while  $\exists k_{j-r} < i' < i$  with  $k_{j-r} < \text{low}(i') = \text{low}(i) \leq k_{j-r+1}$  do
                add column  $i'$  to column  $i$ 
            endwhile
        endfor
    endfor
endfor.

```

The result is the same as that of the Persistence Algorithm in the first section of this chapter; only the order in which the columns are added is different. An easy

connection to persistence arises by considering the monotonic function $f : K \rightarrow \mathbb{R}$ mapping a simplex $\sigma \in K_i - K_{i-1}$ to $f(\sigma) = i$. A leftmost lowest one in ∂_i^j then belongs to a simplex pair of persistence $j - i$. The Spectral Sequence Algorithm thus computes the pairs in the order of non-decreasing index persistence.

Groups and maps. We now interpret the algorithm in terms of groups that make up the spectral sequence of the filtration. Recall the chain groups and boundary maps, $\partial : C_p \rightarrow C_{p-1}$, which form the chain complex defined by K . For each j , we let C_p^j be the group of p -chains of $K_j - K_{j-1}$, and for each chain $c \in C_p^j$, we let $\partial_i^j c$ be the sum of terms of ∂c that lie in $K_i - K_{i-1}$. Suppressing the dimension in the notation for the boundary map, we have $\partial_i^j : C_p^j \rightarrow C_{p-1}^i$ and

$$\partial c = \partial_i^j c + \partial_{j-1}^j c + \dots + \partial_1^j c.$$

The block ∂_i^j in the boundary matrix represents the maps ∂_i^j simultaneously for all dimensions. In spectral sequences, we approximate ∂ by the sum of maps ∂_j^j to ∂_i^j and then decrease i . The spectral sequence itself consists of a collection of groups $E_{p,q}^r$ and maps $d_{p,q}^r$ between them. To describe them, we break with the convention of using p for the dimension. Instead, we follow the convention entrenched in the spectral sequence literature in which the first subscript, p , identifies the block of columns, the sum of subscripts, $p+q$, gives the dimension, and the superscript, r , counts the phases in the iteration.

As usual, we think of the columns of the boundary matrix as generators of the chain groups. Limiting our attention to the p -th block of columns, ∂^p , we get the groups of $(p+q)$ -chains of $K_p - K_{p-1}$, for all q . If we further limit ∂^p to the blocks of rows ∂_i to ∂_p , we effectively ignore any boundary in K_{i-1} . For $i = p$, this is equivalent to taking the relative chain groups, $C_{p+q}(K_p, K_{p-1})$. For $i < p$, we have a subgroup of the relative chain group $C_{p+q}(K_p, K_{i-1})$, namely the one generated by the $(p+q)$ -simplices in $K_p - K_{p-1}$; see Figure VII.13. For what follows, it is

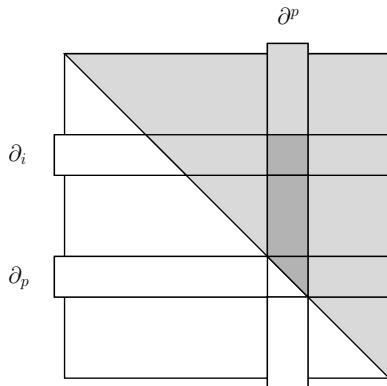


Figure VII.13: The darker shaded portion of the p -th block of columns represents the chains of $K_p - K_{p-1}$ and their boundaries in $K_p - K_{i-1}$.

important to remember that the boundary matrix, ∂ , represents simplices of all dimensions at once. Hence, each block will correspond to a sequence of groups, one for each dimension.

The E^0 -term of the spectral sequence. To prepare for the first phase of the algorithm, we focus on the diagonal blocks of the boundary matrix. Fixing $r = 0$, we write $E_{p,q}^0 = C_{p+q}^p$ for the group of $(p+q)$ -chains of $K_p - K_{p-1}$. Fixing p and varying q , these groups are generated by the p -th block of columns. Furthermore, we let

$$d_{p,q}^0 : E_{p,q}^0 \rightarrow E_{p,q-1}^0$$

be defined by the $(p+q)$ -dimensional boundary map restricted to the block ∂_p^p . In other words, $d_{p,q}^0$ is ∂_p^p applied to $(p+q)$ -chains. We note that $E_{p,q}^0$ is isomorphic to the relative chain group $C_{p+q}(K_p, K_{p-1})$ and $d_{p,q}^0$ agrees with the corresponding relative boundary map. It follows that the maps satisfy the Fundamental Lemma of Homology, that is, $d_{p,q-1}^0 \circ d_{p,q}^0 = 0$. Indeed, a codimension-2 face of a $(p+q)$ -simplex in $K_p - K_{p-1}$ either does not belong to $K_p - K_{p-1}$ or it does, but then both codimension-1 faces that contain it also belong to $K_p - K_{p-1}$. Hence, we get a chain complex,

$$\dots \rightarrow E_{p,q+1}^0 \rightarrow E_{p,q}^0 \rightarrow E_{p,q-1}^0 \rightarrow \dots,$$

in which the maps are implied. It is customary to draw this chain complex vertically, and adding the chain complexes for the other diagonal blocks, we get a 2-dimensional grid of groups, as shown in Figure VII.14. To reduce the clutter, we omit the arrows that connect the groups in each vertical line from top to bottom. We call this the E^0 -term of the spectral sequence, noting that a vertical line in the grid contains all groups represented by a diagonal block of the boundary matrix.

$$\begin{array}{ccccccc} & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & E_{1,1}^0 & E_{2,1}^0 & E_{3,1}^0 & E_{4,1}^0 & E_{5,1}^0 & \cdots \\ \cdots & E_{1,0}^0 & E_{2,0}^0 & E_{3,0}^0 & E_{4,0}^0 & E_{5,0}^0 & \cdots \\ \cdots & E_{1,-1}^0 & E_{2,-1}^0 & E_{3,-1}^0 & E_{4,-1}^0 & E_{5,-1}^0 & \cdots \\ \cdots & 0 & E_{2,-2}^0 & E_{3,-2}^0 & E_{4,-2}^0 & E_{5,-2}^0 & \cdots \\ \cdots & 0 & 0 & E_{3,-3}^0 & E_{4,-3}^0 & E_{5,-3}^0 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \end{array}$$

Figure VII.14: The E^0 -term of the spectral sequence. We have maps going vertically downward, from $E_{p,q}^0$ to $E_{p,q-1}^0$ for every choice of p and q .

The E^1 -term. After interpreting the diagonal blocks of the original boundary matrix in terms of relative chain groups, we now push this interpretation through

the phases of the algorithm. For the first phase, we take the homology of the above vertical complexes and define $E_{p,q}^1 = \ker d_{p,q}^0 / \text{im } d_{p,q+1}^0$. An element of $E_{p,q}^1$ is thus the equivalence class of a chain $c \in C_{p+q}^p$ with $\partial_p^p c = 0$, where two chains are equivalent if their difference lies in the image of ∂_p^p , taking of course the boundary map that applies to chains of one higher dimension. In other words, the element is a relative homology class and more generally $E_{p,q}^1 \simeq H_{p+q}(K_p, K_{p-1})$. Representatives of $E_{p,q}^1$ are computed by reducing the matrix ∂_p^p , which is what the algorithm does in Phase $r = 1$. The zero columns in ∂_p^p correspond to simplices that give birth and represent cycles. Some are paired and have zero persistence since their classes come and go within $K_p - K_{p-1}$. Others are not paired, and their cycles are the generators of $E_{p,q}^1$. Next, we let

$$d_{p,q}^1 : E_{p,q}^1 \rightarrow E_{p-1,q}^1$$

be defined by the $(p+q)$ -th boundary map restricted to ∂_p^{p-1} . Recall that an element in $E_{p,q}^1$ is represented by a relative $(p+q)$ -cycle, c . Hence, $\partial_p^p c = 0$, but $\partial_p^{p-1} c$ is possibly non-zero and represents a class in $E_{p-1,q}^1$. All this sounds complicated, but it is rather straightforward if interpreted in terms of the boundary matrix after one phase of the algorithm. As before, the boundary maps satisfy the Fundamental Lemma of Homology, $d_{p-1,q}^1 \circ d_{p,q}^1 = 0$, so we again get a chain complex:

$$\dots \rightarrow E_{p+1,q}^1 \rightarrow E_{p,q}^1 \rightarrow E_{p-1,q}^1 \rightarrow \dots .$$

Going back to the grid in Figure VII.14, we can see these complexes as horizontal lines going from right to left. Of course, we are now in the next phase, so we need to substitute $r = 1$ for the superscript 0 everywhere. This is the E^1 -term of the spectral sequence.

The E^2 -term. We take one more step before appealing to induction, taking the homology of the horizontal complexes, $E_{p,q}^2 = \ker d_{p,q}^1 / \text{im } d_{p+1,q}^1$. An element of $E_{p,q}^2$ is the equivalence class of the sum of a chain $c \in C_{p+q}^p$ and another chain $c' \in C_{p+q}^{p-1}$. The chains satisfy $\partial_p^p c = 0$ and $\partial_{p-1}^p c + \partial_{p-1}^{p-1} c' = 0$, and being equivalent means that the difference lies in $\text{im } \partial_p^p + \text{im } \partial_{p-1}^p + \text{im } \partial_{p-1}^{p-1}$. The group $E_{p,q}^2$ is not a relative homology group by itself but a subgroup of one, namely $E_{p,q}^2 \oplus E_{p-1,q+1}^1 \simeq H_{p+q}(K_p, K_{p-2})$. Representatives of $E_{p,q}^2$ are computed by reducing the double block of matrices $\partial_p^p, \partial_{p-1}^{p-1}, \partial_p^{p-1}, \partial_{p-1}^p$. The first two have already been reduced, and the third is zero. Phase $r = 2$ completes the reduction of the double block for the remaining fourth matrix. Next, we let

$$d_{p,q}^2 : E_{p,q}^2 \rightarrow E_{p-2,q+1}^2$$

be defined by the $(p+q)$ -th boundary map restricted to ∂_{p-2}^p . By construction, an element of $E_{p,q}^2$ is represented by a $(p+q)$ -chain, c , whose boundary in $K_p - K_{p-2}$ is empty. Its boundary in $K_{p-2} - K_{p-3}$ is possibly non-empty and represents a class in $E_{p-2,q+1}^2$, the image of the class of c in $E_{p,q}^2$. Taking the thus defined boundary map twice gives zero again, so we get a chain complex,

$$\dots \rightarrow E_{p+2,q-1}^2 \rightarrow E_{p,q}^2 \rightarrow E_{p-2,q+1}^2 \rightarrow \dots ,$$

similar to before. Going back to the grid in Figure VII.14, we see this complex along a line of slope one half going from right to left. In other words, the groups are connected by knight moves in chess, two to the left and one up. Of course, we are now in the next phase, so we need to substitute $r = 2$ for the superscript 0 everywhere. This is the E^2 -term of the spectral sequence.

Iteration. The process continues, and for general phase numbers r , the maps take the topologist's chess move, that is, r steps to the left and $r - 1$ steps up:

$$d_{p,q}^r : E_{p,q}^r \rightarrow E_{p-r,q+r-1}^r.$$

This gives a set of chain complexes, and we take homology to enter the next phase. Since K is finite, the maps are eventually zero and the sequence converges to a limit term, $E^r = E^\infty$ for r large enough. The homology groups of K are obtained by taking direct sums along the diagonal lines in the limit term for which the dimension is constant.

Before reaching the limit term, we may consider each class in $E_{p,q}^r$ as generated by an “almost” cycle of dimension $p + q$. This is a chain whose boundary in $K_p - K_{p-r}$ is empty but may have non-empty boundary in K_{p-r} . It is either an essential cycle of K , or a cycle of persistence at least r , assuming the monotonic function $f : K \rightarrow \mathbb{R}$ that maps $\sigma \in K_p - K_{p-1}$ to $f(\sigma) = p$, as before. This leads to the following summary connection between persistence and spectral sequences.

SPECTRAL SEQUENCE THEOREM. The total rank of the groups of dimension $p+q$ after $r \geq 1$ phases of the Spectral Sequence Algorithm equals the number of points in the $(p+q)$ -th persistence diagram of f whose persistence is r or larger; that is,

$$\sum_{p=1}^n \text{rank } E_{p,q}^r = \text{card } \{a \in \text{Dgm}_{p+q}(f) \mid \text{pers}(a) \geq r\},$$

where q decreases as p increases so that the dimension remains constant.

In the limit, for r large enough, we have $\sum_{p=1}^n \text{rank } E_{p,q}^r = \text{rank } H_{p+q}(K)$ equal to the number of points in the $(p+q)$ -th persistence diagram whose persistence is infinite.

Bibliographic notes. A comprehensive account of spectral sequences can be found in [109]. The treatment in this section follows the more concise presentation in the survey of persistent homology [57]. Similar to persistent homology, working over a field is crucial for the construction of spectral sequences. Over \mathbb{Z} , there are extension problems to solve because of torsion; see [24].

Exercises

The credit assignment reflects a subjective assessment of difficulty. A typical question can be answered using knowledge of the material combined with some thought

and analysis.

1. **Tetrahedron complex** (one credit). Let K consist of a tetrahedron and its faces.
 - (i) Apply the matrix reduction algorithm to the filtration of K obtained by adding the simplices in the order of dimension.
 - (ii) Do any of the three diagrams depend on the way you order the simplices of the same dimension?
2. **Matrix reduction revisited** (two credits). Change the standard matrix reduction implementation of the persistence algorithm described in Section VII.1 by adding each j -th column to columns on its right rather than adding columns on its left to it. Specifically, consider the following implementation.

```
R = ∂;
for j = 1 to m do
  while there exists  $j_0 > j$  with  $low(j_0) = low(j)$  do
    add column  $j$  to column  $j_0$ 
  endwhile
endfor.
```

- (i) Show that this implementation of the persistence algorithm generates the same lowest 1s as the standard matrix reduction implementation.
 - (ii) Give an example for which this and the standard implementation of the persistence algorithm compute different reduced matrices.
3. **Sublevel sets** (two credits). Let $f : |K| \rightarrow \mathbb{R}$ be a piecewise linear function defined by its values at the vertices, $f(u_1) < f(u_2) < \dots < f(u_n)$. Let b be strictly between $f(u_i)$ and $f(u_{i+1})$, for some $1 \leq i \leq n - 1$, and recall that the sublevel set defined by b is $f^{-1}(-\infty, b]$.
 - (i) Prove that the sublevel sets defined by b and by $f(u_i)$ have the same homotopy type.
 - (ii) Draw an example for the case in which the sublevel sets defined by b and by $f(u_{i+1})$ have the same homotopy types, and another example for the case in which they have different homotopy types.
4. **Graphs without branching** (three credits). Let K be a 1-dimensional simplicial complex in which each vertex belongs to one or two edges. In other words, K is a simple graph whose components are paths and closed curves. Show that the sparse matrix implementation of the persistence algorithm described in Section VII.2 takes time proportional to the number of simplices in K .
5. **Persistence diagram** (one credit). Draw a genus-3 torus, consider its height function, and draw the non-trivial persistence diagrams of the function. Distinguish between points in the ordinary, extended, and relative subdiagrams.
6. **Breaking symmetry** (two credits). Design a topological space \mathbb{X} and a continuous function $f : \mathbb{X} \rightarrow \mathbb{R}$ such that

- (i) the persistence diagrams violate the Persistence Duality Theorem in Section VII.3;
 - (ii) the persistence diagrams violate the Persistence Symmetry Theorem in the same section.
7. **Matrix reduction once again** (one credit). Prove that the reduced matrix computed by the spectral sequence algorithm in Section VII.4 is the same as that generated by the persistence algorithm in Section VII.1.
8. **Parallel matrix reduction** (three credits). First, rewrite the Spectral Sequence Algorithm of Section VII.4 for the case in which each block, $K_j - K_{j-1}$, consists of a single simplex. Second, show that the thus simplified algorithm can be run on a parallel computer architecture using n processors taking time at most proportional to n^2 .

Chapter VIII

Stability

Persistence is a measure-theoretic concept built on top of algebraic structures. Its most important property is the stability under perturbations of the data. In other words, small changes in the data imply at most small changes in the measured persistence. This has major ramifications, including the study of 1-parameter families and the comparison and classification of shapes. Of particular importance are biological shapes, with their sheer endless variety in the midst of unmistakable similarity and delicate variation. This book touches upon this fascinating topic, and we foresee future inroads based on the notion of persistent homology as developed here.

VIII.1 1-parameter Families

In this section, we study how continuous change of the data affects the measured persistence. We focus on the structural effects and their computation. A consequence of the analysis is a first proof of stability.

Straight-line homotopy. Let $f : K \rightarrow \mathbb{R}$ and $g : K \rightarrow \mathbb{R}$ be two monotonic functions on the same simplicial complex. We recall that this means that the functions are non-decreasing along increasing chains of the face relation. We use the straight-line homotopy $F : K \times [0, 1] \rightarrow \mathbb{R}$ defined by

$$F(\sigma, t) = (1 - t)f(\sigma) + tg(\sigma)$$

to interpolate between f and g . Define $f_t(\sigma) = F(\sigma, t)$ and note that $f_0 = f$ and $f_1 = g$, as intended. Furthermore, f_t is monotonic for each $t \in [0, 1]$. Indeed, if σ is a face of τ , then $f(\sigma) \leq f(\tau)$ and $g(\sigma) \leq g(\tau)$ and therefore $f_t(\sigma) \leq f_t(\tau)$ for every $t \in [0, 1]$. Hence, we can find a compatible ordering of the simplices, that is, a total order that extends the partial orders defined by f_t and by the face relation. Using this compatible ordering, we compute the persistence diagrams of f_t as explained

in the previous chapter. However, if we somehow already have the diagrams for f , then we may consider modifying them to get the diagrams for f_t . This turns out to be more efficient than recomputing the diagrams provided the two total orders are not too different. To describe exactly what this means, we plot the function values with time, giving us a straight line for each simplex; see Figure VIII.1. It

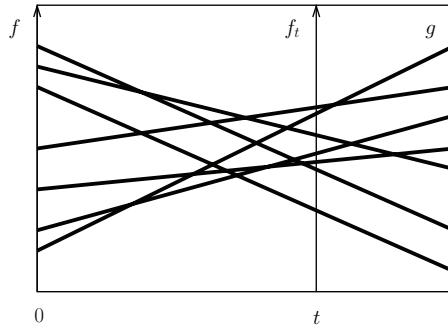


Figure VIII.1: Each line tracks the function value of a simplex as t increases. At any moment $t \in [0, 1]$, we get f_t by intersection with the corresponding vertical line.

is convenient to assume that f and g are injective because this implies that f_t is injective except at finitely many moments t when two or more of the lines cross. To further simplify the situation, we make the generic assumption that no two different pairs of lines cross at the same moment. Equivalently, every f_t has at most one violation of injectivity, namely at most two simplices with the same function value. As we sweep from left to right, in the direction of increasing t , we pass through each violation by transposing the two simplices in the compatible ordering. This motivates us to study the impact of a transposition on persistence.

Matrix decomposition. We recall that we compute the persistence diagrams of $f : K \rightarrow \mathbb{R}$ by reducing the boundary matrix whose rows and columns are ordered like the simplices in a compatible ordering. Starting with $R = \partial$, we perform left-to-right column additions until R is reduced, that is, each non-zero column has its lowest 1 in a unique row. In other words, the mapping from non-zero columns to rows defined by *low* is injective. Each lowest 1 gives a pair of simplices, namely (σ_i, σ_j) if $i = \text{low}(j)$, and a finite off-diagonal point in the p -th persistence diagram, namely $(f(\sigma_i), f(\sigma_j))$ in $\text{Dgm}_p(f)$ with $p = \dim \sigma_i$. It will be convenient to assume a bijection between the lowest 1s and the off-diagonal points in the persistence diagrams. In other words, we assume there are no off-diagonal points at infinity or, equivalently, that every zero column in the reduced matrix corresponds to a row with a lowest 1. We get this property in reduced homology iff K is homologically trivial. This is no loss of generality since we can always add simplices at the end so that they do not alter the earlier homological evolution along the filtration. For example, we can form the cone over a given simplicial complex, which is necessarily homologically trivial.

The reduced matrix can be written as $R = \partial V$, where V keeps track of the column operations. Its j -th column stores the chain whose boundary is stored in the j -th column of R . Since we only use left-to-right column additions, V is upper triangular, with $V[i, i] = 1$ for each i . The matrix V is therefore invertible. Let U be the right inverse of V and note that it is again upper triangular and invertible. Multiplying on the right, we get $RU = \partial VU$ and therefore

$$\partial = RU.$$

We call this an *ru-decomposition* of the boundary matrix, using lowercase letters to emphasize that we mean the form rather than the names of the matrices. Implicit in this definition are the requirements that U be upper triangular and invertible and that R be reduced. We get these properties from the way we compute the matrices, but there are other *ru*-decompositions that may be obtained by other, similar algorithms. Indeed, the *ru*-decomposition of ∂ is not unique, but as noted in the previous chapter, the lowest 1s in the reduced matrix are. The specific question we now ask is how we can update the *ru*-decomposition of the boundary matrix if we transpose two simplices in contiguous positions along the compatible ordering.

Updating the decomposition. Suppose ∂ is the boundary matrix for the ordering of the simplices as $\sigma_1, \sigma_2, \dots, \sigma_m$. We write ∂' for the boundary matrix after transposing σ_i with σ_{i+1} . Letting $P = P_i^{i+1}$ be the corresponding permutation matrix, we have $\partial' = P\partial P$. The difference between P and the unit matrix, I , is localized to the 2-by-2 submatrix for which $P[i, i] = P[i+1, i+1] = 0$ and $P[i, i+1] = P[i+1, i] = 1$. Multiplying by P from the left exchanges the two rows, and multiplying by P from the right exchanges the two columns. Note also that P is its own inverse, that is, $PP = I$. We therefore get

$$\partial' = P\partial P = PRUP = (PRP)(PUP).$$

But this is not necessarily an *ru*-decomposition of the new boundary matrix. It fails to be one if $R' = PRP$ is not reduced or if $U' = PUP$ is not upper triangular. We will now show that either deficiency can be remedied with little effort, namely a constant number of row and column operations.

The only way R' can fail to be reduced is when rows i and $i+1$ of R both contain a lowest 1, $i = \text{low}(k)$ and $i+1 = \text{low}(l)$, and row i has a 1 in column l as well. Notice that $i, i+1 < k, l$. There are two cases, distinguished by $k < l$ and $l < k$. In both cases, we add the left column to the right column before we do the transposition. This fixes the deficiency, as illustrated in Figure VIII.2, by changing R before we even make the transposition of σ_i with σ_{i+1} .

The only way U' can fail to be upper triangular is if $U[i, i+1] = 1$. We fix this deficiency by adding row $i+1$ to row i in U and adding column i to column $i+1$ in R . Letting $S = S_i^{i+1}$ be the matrix whose only difference from the identity matrix is that $S[i, i+1] = 1$, we thus consider SU and RS . Since $PP = I$ and $SS = I$, this does not change the matrix product; that is, $\partial' = (PRSP)(PSUP) = PRUP$, as

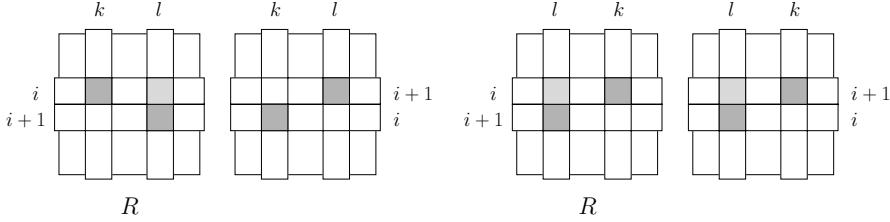


Figure VIII.2: After swapping rows i and $i + 1$ in R , the matrix would be no longer reduced. We thus add the left to the right column before exchanging the two rows.

before. With this modification, $PSUP$ is upper triangular, but $PRSP$ may again fail to be reduced. If column i is zero or $low(i) < low(i + 1)$, then multiplying by S preserves the lowest 1s and RS is reduced. In this case, we have an *ru*-decomposition after the transposition. On the other hand, if column $i + 1$ is zero while column i is not or if $low(i) > low(i + 1)$, as in Figure VIII.3, then we need to make the lowest 1s unique again. To do this, we add column $i + 1$ to column i after the transposition resulting in a left-to-right column addition. This repairs all deficiencies, and we have an *ru*-decomposition of ∂' .

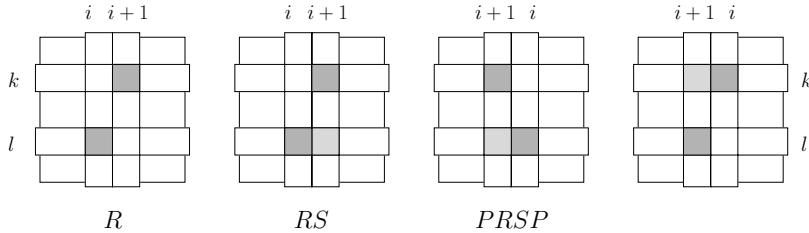


Figure VIII.3: After adding column i to $i + 1$ and exchanging the two columns, the matrix $PRSP$ is no longer reduced. Adding column $i + 1$ to i after the transposition finally produces a reduced matrix.

How the pairing changes. It takes additional work to understand which transpositions have an effect on the pairing. The ones that do we call *switches*. Recall that each lowest 1 establishes a correspondence between a positive simplex (a row) and a negative simplex (a column). For example, in Figure VIII.2 on the left, we have the pairs (σ_i, σ_k) and (σ_{i+1}, σ_l) , which are preserved throughout the transposition. On the right, we have the same two pairs but they change to (σ_i, σ_l) and (σ_{i+1}, σ_k) . This identifies the transposition as a switch. In Figure VIII.3, we have the pairs (σ_k, σ_{i+1}) and (σ_l, σ_i) , which change to (σ_k, σ_i) and (σ_l, σ_{i+1}) , again a switch.

As a rule of thumb, most transpositions are not switches. For example, if σ_i and σ_{i+1} do not have the same dimension, then their transposition does not require any changes other than the obligatory swapping of rows and columns. Even if they have

the same dimension but if σ_i is positive and σ_{i+1} is negative, then the transposition cannot be a switch. This is because row i has no lowest 1, so $R' = PRP$ is reduced and requires no further effort. Similarly, column i of R is zero, so we can set $U[i, i+1] = 0$ to make sure $U' = PUP$ is upper triangular, if necessary. In words, the ru -decomposition is maintained without any of the repair operations that change the pairing. However, the remaining three combinations of types can be switches, and we see an example of each in Figure VIII.4. We get a switch between two

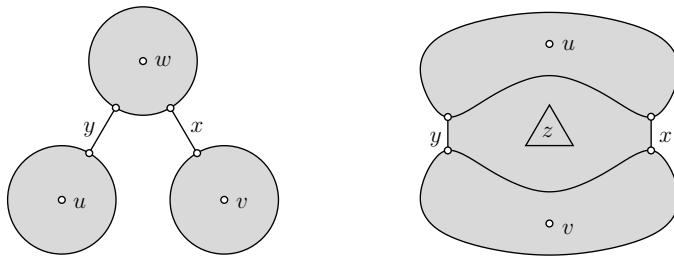


Figure VIII.4: The vertices u, v, w are the oldest in their respective components, which are eventually joined by the edges x and y . On the right, the two edges form a hole, which is eventually filled by the triangle z .

positive vertices, v and w , when we go from $uvwxy$ to $uwvxy$ in the drawing on the left. Indeed, the pairs (v, y) and (w, x) before the transposition of v and w change to (v, x) and (w, y) after the transposition. We get a switch between two negative edges, x and y , when we go from $uvwxy$ to $uvwyx$, again on the left. Indeed, the transposition of x and y produces the same change between pairs as in the previous example. Finally, we get a switch between a negative edge, x , and a positive edge, y , when we go from $uvxyz$ to $uvyxz$ in the drawing on the right. Indeed, the pairs (v, x) and (y, z) before the transposition of x and y change to (v, y) and (x, z) after the transposition. The last switch is the most interesting of all. Besides changing the pairing, it convinces the negative x to become positive and the positive y to become negative. The two edges thus contribute to different persistence diagrams before and after the transposition.

Summary. When we transpose σ_i and σ_{i+1} , we touch only the columns of σ_i and σ_{i+1} and of the simplices σ_k and σ_l paired with them. The changes are therefore limited to these two pairs. Furthermore, there are no changes unless the transposed simplices have the same dimension. Assuming $p = \dim \sigma_i = \dim \sigma_{i+1}$, the other two simplices have dimension $p - 1$ and $p + 1$. The only possible change is therefore that the transposed simplices trade places. We state this result for later reference.

TRANSPOSITION LEMMA. Let ∂ and ∂' be the boundary matrices for compatible orderings of two monotonic functions on a simplicial complex that differ by a single transposition of two contiguous simplices, σ_i and σ_{i+1} . Then the pairings defined by ru -decompositions $\partial = RU$ and $\partial' = R'U'$ differ only if $\dim \sigma_i = \dim \sigma_{i+1}$, and if they differ, then only by σ_i and σ_{i+1} trading places.

The computational effort for updating the *ru*-decomposition is small, namely a constant number of row and column operations, each computable in time proportional to the number of simplices. Returning to our two monotonic functions, $f, g : K \rightarrow \mathbb{R}$, we have m simplices and thus at most $\binom{m}{2}$ transpositions to go from a compatible ordering for f to a compatible ordering for g . To get started, we compute the persistence diagrams of f in m^3 time using the algorithm explained in Section VII.1. Thereafter, we spend m time per transposition and therefore $m\binom{m}{2} < m^3$ time in total until we arrive at the persistence diagrams of g . This is roughly the same amount of time required to compute the diagrams of g from scratch, at least in the worst case. However, going through the transposition has the advantage that we get the interpolating diagrams for free.

Bibliographic notes. The material of this section is fashioned after [38], where continuous families of persistence diagrams are proposed as a tool to study parametrized families of functions. As explained, the algorithm constructs these diagrams by maintaining the *ru*-decomposition of the boundary matrix through a sequence of transpositions scheduled by sweeping an arrangement of lines. We can find these transpositions in logarithmic time each by sorting the crossings or in constant time each by sweeping the arrangement topologically [56].

VIII.2 Stability Theorems

Like any good measurement device, persistence gives similar readings for similar functions. We make this statement precise for two notions of similarity between persistence diagrams. The bottleneck distance is the cruder of the two but leads to a more general result. The Wasserstein distance is more sensitive to details in the diagrams but requires additional properties to be stable.

Bottleneck distance. Recall that a persistence diagram is a multiset of points in the extended plane, $\bar{\mathbb{R}}^2$. Under the assumptions on the input functions considered in this book, the diagram consists of finitely many points above the diagonal. To this finite multiset, we add the infinitely many points on the diagonal, each with infinite multiplicity. These extra points are not essential to the diagram, but their presence simplifies upcoming definitions and results. Now let X and Y be two persistence diagrams. To define the distance between them, we consider bijections $\eta : X \rightarrow Y$ and record the supremum of the distances between corresponding points for each. Measuring distance between points $x = (x_1, x_2)$ and $y = (y_1, y_2)$ with L_∞ -norm $\|x - y\|_\infty = \max\{|x_1 - y_1|, |x_2 - y_2|\}$ and taking the infimum over all bijections, we get the *bottleneck distance* between the diagrams:

$$W_\infty(X, Y) = \inf_{\eta: X \rightarrow Y} \sup_{x \in X} \|x - \eta(x)\|_\infty.$$

As illustrated in Figure VIII.5, we can draw squares of side length twice the bottleneck distance centered at the points of X so that each square contains the

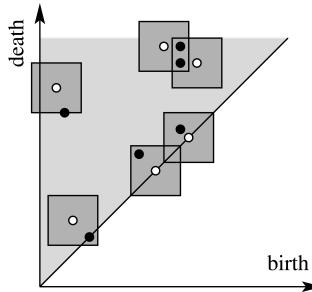


Figure VIII.5: The superposition of two persistence diagrams consisting of the white and the black points. Only the marked points on the diagonal correspond to off-diagonal points in the other diagram. The bottleneck distance is half the side length of the squares illustrating the bijection.

corresponding point of Y . Clearly, $W_\infty(X, Y) = 0$ iff $X = Y$. Furthermore, $W_\infty(X, Y) = W_\infty(Y, X)$ and $W_\infty(X, Z) \leq W_\infty(X, Y) + W_\infty(Y, Z)$. We see that W_∞ satisfies all axioms of a metric and thus deserves to be called a distance.

Bottleneck stability. Letting $f, g : K \rightarrow \mathbb{R}$ be two monotonic functions, we consider the straight-line homotopy $f_t = (1-t)f + tg$, as in the previous section. This gives a monotonic function f_t with a persistence diagram for each dimension p and each $t \in [0, 1]$. Fixing a dimension p , the family of persistence diagrams is a multiset in $\mathbb{R}^2 \times [0, 1]$. Drawing t along a third coordinate axis, we get a 3-dimensional visualization of how the persistent homology evolves as we go from $f_0 = f$ to $f_1 = g$. To describe this, we assume that K has no non-trivial (reduced) homology, as in the previous section. Adding the third coordinate, each off-diagonal point of $X_t = \text{Dgm}_p(f_t)$ is of the form $x(t) = (f_t(\sigma), f_t(\tau), t)$, where σ and τ are simplices in K . The point represents the fact that when we construct K by adding the simplices in the order defined by f_t , then adding σ gives birth to a p -dimensional homology class and adding τ gives death to the same. There are only finitely many values at which the pairing of the simplices changes, and we denote these as $0 = t_0 < t_1 < \dots < t_n < t_{n+1} = 1$. Within each interval (t_i, t_{i+1}) , the pairing is constant and each pair σ, τ gives rise to a line segment of points $x(t)$ connecting points in the planes $t = t_i$ and $t = t_{i+1}$. If the endpoint is an off-diagonal point at t_{i+1} , then there is some other unique line segment that begins at that point. This line segment may correspond to the same simplex pair and thus continue on the same straight line, or it may correspond to a different pair created in a switch and make a turn at the shared point. It is also possible that the endpoint lies on the diagonal at t_{i+1} , in which case there is no continuation. In summary, the line segments form polygonal paths that monotonically increase in t . Each path begins at an off-diagonal point in $X = X_0$ or at a diagonal point in some X_{t_i} and ends at an off-diagonal point in $Y = X_1$ or at a diagonal point in some X_{t_j} . We call each polygonal path a *vine* and the multiset of vines a *vineyard*; see Figure VIII.6.

The fact that the points in the family of persistence diagrams form connected

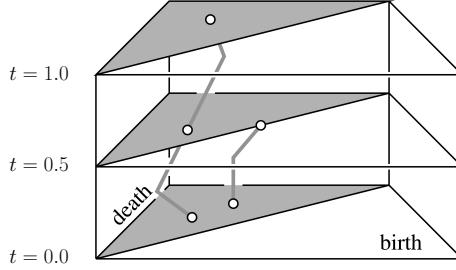


Figure VIII.6: The 1-parameter family of persistence diagrams of the straight-line homotopy between $f_0 = f$ and $f_1 = g$. One point traces out a vine spanning the entire interval while the other merges into the diagonal halfway through the homotopy.

vines is important. It is a way of saying that the persistence diagram is stable. To further quantify this notion, we differentiate $x(t) = (1-t)(f(\sigma), f(\tau), 0) + t(g(\sigma), g(\tau), 1)$ to get $\frac{\partial x}{\partial t}(t) = (g(\sigma) - f(\sigma), g(\tau) - f(\tau), 1)$. Projecting the endpoints of the line segment back into \mathbb{R}^2 , we get two points whose L_∞ -distance is $t_{i+1} - t_i$ times the larger of the differences between f and g at the two simplices. Letting v be the simplex in K that maximizes this difference, we get the L_∞ -distance between the two functions, $\|f - g\|_\infty = |f(v) - g(v)|$. This is also an upper bound on the slope of any line segment in the vineyard and therefore an upper bound on the L_∞ -distance between the projected endpoints of any vine.

STABILITY THEOREM FOR FILTRATIONS. Let K be a simplicial complex and $f, g : K \rightarrow \mathbb{R}$ two monotonic functions. For each dimension p , the bottleneck distance between the diagrams $X = \text{Dgm}_p(f)$ and $Y = \text{Dgm}_p(g)$ is bounded from above by the L_∞ -distance between the functions, $W_\infty(X, Y) \leq \|f - g\|_\infty$.

Tame functions. To apply the Stability Theorem, it is convenient to get it into a form that allows for more general functions. According to the Simplicial Approximation Theorem in Chapter III, every continuous function on a triangulable topological space can be approximated by a piecewise linear function, and as shown in Chapter VII, for every piecewise linear function there is a monotonic function that generates the same persistence diagrams. It is therefore not surprising that what we said about filtrations can indeed be generalized. We explain this for functions that satisfy a mild tameness condition.

Let \mathbb{X} be triangulable and $f : \mathbb{X} \rightarrow \mathbb{R}$ continuous. Given a threshold $a \in \mathbb{R}$, recall that the sublevel set consists of all points $x \in \mathbb{R}$ with function value less than or equal to a , $\mathbb{X}_a = f^{-1}(-\infty, a]$. Similar to the complexes in a filtration, the sublevel sets are nested and give rise to a sequence of homology groups connected by maps induced by inclusion, one for each dimension. Writing $f_p^{a,b} : \mathbb{H}_p(\mathbb{X}_a) \rightarrow \mathbb{H}_p(\mathbb{X}_b)$ for the map from the p -th homology group of the sublevel set at a to that at b , we call its image a *persistent homology group*, as before. The corresponding *persistent Betti number* is $\beta_p^{a,b} = \text{rank im } f_p^{a,b}$. As long as the topology of the sublevel set

does not change, the maps between the homology groups are isomorphisms. We thus call $a \in \mathbb{R}$ a *homological critical value* if there is no $\varepsilon > 0$ for which $f_p^{a-\varepsilon, a+\varepsilon}$ is an isomorphism for each dimension p . Finally, we call f *tame* if it has only finitely many homological critical values and all homology groups of all sublevel sets have finite rank. The main motivation for this definition is the relative ease with which we can define persistence diagrams. Letting $a_1 < a_2 < \dots < a_n$ be the homological critical values of f , we choose interleaved values b_0 to b_n with $b_{i-1} < a_i < b_i$ for all i . Adding $b_{-1} = a_0 = -\infty$ and $a_{n+1} = b_{n+1} = \infty$, we then consider the corresponding sequence of homology groups,

$$0 = H_p(\mathbb{X}_{b_{-1}}) \rightarrow H_p(\mathbb{X}_{b_0}) \rightarrow \dots \rightarrow H_p(\mathbb{X}_{b_n}) \rightarrow H_p(\mathbb{X}_{b_{n+1}}) = H_p(\mathbb{X}),$$

and the maps between them. For $0 \leq i < j \leq n + 1$, the *multiplicity* of the pair a_i, a_j is now defined as $\mu_p^{a_i, a_j} = (\beta_p^{b_i, b_{j-1}} - \beta_p^{b_i, b_j}) - (\beta_p^{b_{i-1}, b_{j-1}} - \beta_p^{b_{i-1}, b_j})$. To get the p -th *persistence diagram* of f , we draw each point (a_i, a_j) with multiplicity $\mu_p^{a_i, a_j}$, and we add the points of the diagonal, each with infinite multiplicity. With these definitions, we have the following stability result, which we state without proof, illustrating it in Figure VIII.7.

STABILITY THEOREM FOR TAME FUNCTIONS. Let \mathbb{X} be a triangulable topological space and let $f, g : \mathbb{X} \rightarrow \mathbb{R}$ be two tame functions. For each dimension p , the bottleneck distance between $X = \text{Dgm}_p(f)$ and $Y = \text{Dgm}_p(g)$ is bounded by the L_∞ -distance between the functions, $W_\infty(X, Y) \leq \|f - g\|_\infty$.

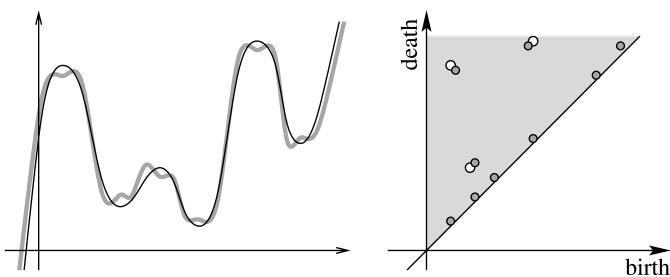


Figure VIII.7: Left: two functions with small L_∞ -distance. Right: the corresponding two persistence diagrams with small bottleneck distance.

Wasserstein distance. A drawback of the bottleneck distance is its insensitivity to details of the bijection beyond the furthest pair of corresponding points. To remedy this shortcoming, we introduce the *degree- q Wasserstein distance* between X and Y for any positive real number q . It takes the sum of q -th powers of the L_∞ -distances between corresponding points, again minimizing over all bijections:

$$W_q(X, Y) = \left[\inf_{\eta: X \rightarrow Y} \sum_{x \in X} \|x - \eta(x)\|_\infty^q \right]^{1/q}.$$

As suggested by our notation, the bottleneck distance is the limit of the Wasserstein distance as q goes to infinity. Similar to the bottleneck distance, it is straightforward to verify that W_q satisfies the requirements of a metric and thus deserves to be called a distance.

It should be obvious that we cannot substitute the degree- q Wasserstein distance for the bottleneck distance and expect that the Stability Theorem for Tame Functions still holds. Indeed, we can approximate a function $f : \mathbb{R} \rightarrow \mathbb{R}$ with a function g that has arbitrarily many wrinkles without deviating from f by more than some positive ε ; see Figure VIII.7. Each wrinkle generates a point with persistence about 2ε in the 0-th persistence diagram. Making the wrinkles narrow, we can get an arbitrarily large number and therefore an arbitrarily large Wasserstein distance between the diagrams of f and g .

Wasserstein stability. Although a general stability result like the one for the bottleneck distance is out of reach, we get stability under the Wasserstein distance for a reasonably large class of functions. Let \mathbb{X} be a metric space, that is, a topological space for which the distance between points $x, y \in \mathbb{X}$, denoted as $\|x - y\|$, is well defined. A function $f : \mathbb{X} \rightarrow \mathbb{R}$ is *Lipschitz* if there is a constant C such that $|f(x) - f(y)| \leq \|x - y\|$ for all points $x, y \in \mathbb{X}$. Without loss of generality, we only consider Lipschitz functions with constant $C = 1$. This condition prevents narrow wrinkles. Indeed, each wrinkle now requires an amount of space that relates to its persistence. It is therefore not possible to crowd arbitrarily many wrinkles together without shrinking their persistence. What we suggest here is a packing argument, the metric version of the combinatorial pigeonhole principle, but homology classes can interact so that the packing argument cannot be applied directly. Indeed, making it a rigorous proof is work which we would rather skip. Instead, we introduce the precise conditions on the space \mathbb{X} for which we can prove stability of persistence.

Assume \mathbb{X} is triangulable and consider a triangulation, that is, a simplicial complex K together with a homeomorphism $\phi : |K| \rightarrow \mathbb{X}$. Letting its mesh be the maximum distance between the images of two points of the same simplex in K , we define $N(r)$ as the minimum number of simplices in a triangulation with mesh at most r . We say the triangulations of \mathbb{X} *grow polynomially* if there are constants c and j such that $N(r) \leq c/r^j$. Finally, we define the *degree- k total persistence* of a persistence diagram X as the sum of k -th powers of the persistences of its points, $\Phi^k(X) = \sum_{x \in X} \text{pers}(x)^k$. To finesse the difficulties caused by points with infinite persistence, we restrict the sum to the finite points in X . The main technical insight is that polynomial growth implies bounded total persistence. Specifically, if \mathbb{X} is a metric space whose triangulations grow polynomially with constant exponent j , $f : \mathbb{X} \rightarrow \mathbb{R}$ is Lipschitz, and $X = \text{Dgm}_p(f)$, then $\Phi^k(X)$ is bounded from above by a constant for every $k > j$. The proof of this implication is omitted. For example the d -dimensional sphere is triangulable, and its triangulations grow polynomially, with constant exponent $j = d$. It follows that for every $k > d$, the degree- k total persistence of a Lipschitz function on the sphere is bounded by a constant. Using these ingredients, we are now ready to prove an upper bound on the Wasserstein distance that implies stability for $q > k$.

STABILITY THEOREM FOR LIPSCHITZ FUNCTIONS. Let $f, g : \mathbb{X} \rightarrow \mathbb{R}$ be tame Lipschitz functions on a metric space whose triangulations grow polynomially with constant exponent j . Then there are constants C and $k > j$ no smaller than 1 such that the degree- q Wasserstein distance between $X = \text{Dgm}_p(f)$ and $Y = \text{Dgm}_p(g)$ is $W_q(X, Y) \leq C \cdot \|f - g\|_\infty^{1-k/q}$ for every $q \geq k$.

PROOF. Let $\eta : X \rightarrow Y$ be a bijection that realizes the bottleneck distance; that is, $\|x - \eta(x)\|_\infty \leq \varepsilon = \|f - g\|_\infty$ for each point $x \in X$. In addition, we require that $\|x - \eta(x)\|_\infty \leq \frac{1}{2}[\text{pers}(x) + \text{pers}(\eta(x))]$. Indeed, if this inequality does not hold, then $\text{pers}(x) \leq 2\varepsilon$ and $\text{pers}(\eta(x)) \leq 2\varepsilon$ and we can change the bijection by matching both with points on the diagonal within L_∞ -distance ε . The q -th power of the degree- q Wasserstein distance is therefore

$$\begin{aligned} W_q(X, Y)^q &\leq \sum_{x \in X} \|x - \eta(x)\|_\infty^q \\ &\leq \varepsilon^{q-k} \sum_{x \in X} \|x - \eta(x)\|_\infty^k \\ &\leq \frac{\varepsilon^{q-k}}{2^k} \sum_{x \in X} [\text{pers}(x) + \text{pers}(\eta(x))]^k \\ &\leq \frac{\varepsilon^{q-k}}{2^k} \sum_{x \in X} [(2\text{pers}(x))^k + (2\text{pers}(\eta(x)))^k], \end{aligned}$$

where the last step uses the fact that taking the k -th power is convex. The sum is 2^k times the degree- k total persistence of X plus that of Y , which gives $W_q(X, Y)^q \leq \varepsilon^{q-k}[\Phi^k(X) + \Phi^k(Y)]$. By assumption, the degree- k total persistence is bounded by a constant. Taking the q -th root thus gives the claimed inequality. \square

Bibliographic notes. Vineyards have been introduced as a tool for studying parametric families of functions in [38]. The proof that the vines in it are connected paths is equivalent to establishing the stability for monotonic functions under the bottleneck distance between diagrams. The first proof of stability goes back to Cohen-Steiner, Edelsbrunner, and Harer [34] who used a homological algebra argument to establish it for tame functions. A further generalization to 1-parameter families of vector spaces can be found in [31]. A proof of the Stability Theorem for Lipschitz Functions along with applications in systems biology can be found in [36]. The Wasserstein distance is named after the author of [153]. It is related to optimal transportation as studied by Monge [112] and Kantorovich [89]; see also [148].

VIII.3 Length of a Curve

In this section, we use the stability of persistence to generalize a classic result on curves, proving an inequality connecting the lengths and total curvatures of two curves. At this time, no other proof of this connection is known.

Closed curves. We consider a closed curve $\gamma : \mathbb{S}^1 \rightarrow \mathbb{R}^2$, with or without self-intersections. Assuming γ is smooth, we have derivatives of all orders. The *speed at a point* $\gamma(s)$ is the length of the velocity vector, $\|\dot{\gamma}(s)\|$. We can use it to compute the length as the integral over the curve:

$$\text{length}(\gamma) = \int_{s \in \mathbb{S}^1} \|\dot{\gamma}(s)\| ds.$$

It is convenient to assume a constant speed parametrization, that is, speed $= \|\dot{\gamma}(s)\| = \text{length}(\gamma)/2\pi$ for all $s \in \mathbb{S}^1$. With this assumption, the *curvature at a point* $\gamma(s)$ is the norm of the second derivative divided by the square of the speed, $\kappa(s) = \|\ddot{\gamma}(s)\|/\text{speed}^2$. The reciprocal of the curvature is the radius of the circle that best approximates the shape of the curve at the point $\gamma(s)$. To interpret this formula geometrically, we follow the velocity vector as we trace out the curve. Since its length is constant, it sweeps out a circle of radius speed, as illustrated in Figure VIII.8. The curvature is the speed at which the unit tangent vector sweeps out the

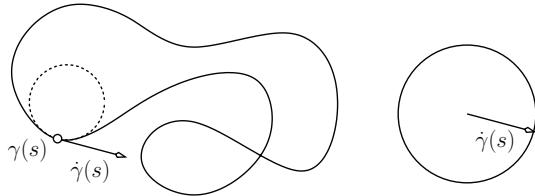


Figure VIII.8: A curve with constant speed parametrization and its velocity vector sweeping out a circle with radius equal to the speed.

unit circle as we move the point with unit speed along the curve. This explains why we divide by the speed twice, first to compensate for the length of the velocity vector and second for the actual speed. The *total curvature* is the distance traveled by the unit tangent vector:

$$\text{curv}(\gamma) = \text{speed} \int_{s \in \mathbb{S}^1} \kappa(s) ds.$$

As an example consider the constant speed parametrization of the circle with radius r , $\gamma(s) = rs$. Writing a point in terms of its angle, we get

$$s = \begin{bmatrix} \cos \varphi \\ \sin \varphi \end{bmatrix}, \quad \gamma(s) = \begin{bmatrix} r \cos \varphi \\ r \sin \varphi \end{bmatrix}, \quad \dot{\gamma}(s) = \begin{bmatrix} -r \sin \varphi \\ r \cos \varphi \end{bmatrix}.$$

We thus have speed $= r$ and $\text{length}(\gamma) = \int \text{speed} ds = 2\pi r$. The curvature is $\kappa(s) = \|\ddot{\gamma}(s)\|/\text{speed}^2 = 1/r$, which is of course independent of the location on the circle. The total curvature is $\text{curv}(\gamma) = \int \frac{r}{r} ds = 2\pi$, which is independent of the radius. Indeed, the unit tangent vector travels once around the unit circle, no matter how small or how big the parametrized circle is.

Integral geometry. The length and total curvature of a curve can also be expressed in terms of integrals of elementary quantities. We begin with the length. Take a unit length line segment in the plane. The lines that cross the line segment at an angle φ form a strip of width $\sin \varphi$. Integrating over all angles gives $\int_{\varphi=0}^{\pi} \sin \varphi d\varphi = [-\cos \varphi]_0^\pi = 2$. In words, the integral of the number of intersections over all lines in the plane is twice the length of the line segment. Since we can approximate a curve by a polygon whose total length approaches that of the curve, the same holds for our curve γ . To express this result formally, we introduce $g_u : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by $g_u(x) = \langle u, x \rangle$, mapping each point $x \in \mathbb{R}^2$ to its height in the direction $u \in \mathbb{S}^1$. The preimage of a value $z \in \mathbb{R}$, $g_u^{-1}(z)$, is the line with normal direction u and offset z . The composition with the curve, $f_u = g_u \circ \gamma$, maps each $s \in \mathbb{S}^1$ to the height of the point $\gamma(s)$. The preimage of this function thus corresponds to points at which the line intersects the curve. We are now ready to formulate the length of the curve in terms of the number of intersections.

CAUCHY-CROFTON FORMULA. The length of a curve in the plane is one quarter of the integral of the number of intersections with lines:

$$\text{length}(\gamma) = \frac{1}{4} \int_{u \in \mathbb{S}^1} \int_{z \in \mathbb{R}} \text{card}(f_u^{-1}(z)) dz du.$$

Here we divide by two twice, once because $\int \sin \varphi d\varphi = 2$ and again because we integrate over all $u \in \mathbb{S}^1$ and therefore over all lines twice. To get an integral geometry expression of the total curvature, we again consider a direction $u \in \mathbb{S}^1$ and the height of the curve in that direction, $f_u : \mathbb{S}^1 \rightarrow \mathbb{R}$. For generic directions u , this height function has a finite number of minima and maxima, as illustrated in Figure VIII.9. Recall that the total curvature is the length traveled by the unit

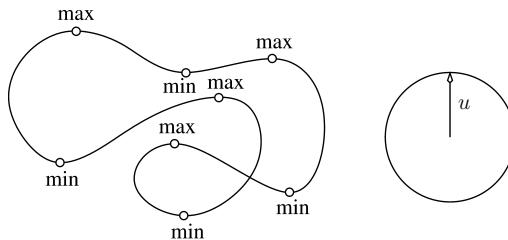


Figure VIII.9: The vertical height function defined on the curve has four local minima which alternate with the four local maxima along the curve.

tangent vector. Equivalently, it is the length traveled by the outward unit normal vector. The number of maxima of f_u is the number of times the unit normal passes $u \in \mathbb{S}^1$, and the number of minima is the number of times it passes $-u \in \mathbb{S}^1$. Writing $\#\text{crit}(f_u)$ for the number of minima and maxima, we get the total curvature by integration.

TOTAL CURVATURE FORMULA. The total curvature of a smooth curve in the plane is half the integral of the number of critical points over all directions:

$$\text{curv}(\gamma) = \frac{1}{2} \int_{u \in \mathbb{S}^1} \#\text{crit}(f_u) du.$$

The integral in the above formula can be interpreted as 2π times the average number of critical points, where the average is taken over all directions. Hence the total curvature is π times this average.

Theorems relating length with total curvature. Suppose the image of γ fits inside the unit disk in the plane, $\text{im } \gamma \subseteq \mathbb{B}^2$. Then γ must turn to avoid crossing the boundary circle of the disk. We can therefore expect that the total curvature is bounded from below by some constant times the length. A classic result in geometry asserts that this constant is one.

FÁRY THEOREM. Let $\gamma : \mathbb{S}^1 \rightarrow \mathbb{R}$ be a smooth closed curve with $\text{im } \gamma \subseteq \mathbb{B}^2$. Then its length is at most its total curvature, $\text{length}(\gamma) \leq \text{curv}(\gamma)$.

To generalize this result, we consider two curves, $\gamma, \gamma_0 : \mathbb{S}^1 \rightarrow \mathbb{R}^2$, and the ‘shortest leash distance’ between them. Specifically, we trace out both curves simultaneously and connect the two moving points by a leash so that their distance can never exceed the length of that leash. Formally, this concept is known as the *Fréchet distance* between the curves. To define it, we record the leash length for a homeomorphism $\eta : \mathbb{S}^1 \rightarrow \mathbb{S}^1$ and take the infimum over all homeomorphisms, $F(\gamma, \gamma_0) = \inf_{\eta} \max_s \|\gamma(s) - \gamma_0(\eta(s))\|$. This notion of distance does not depend on the parametrizations of the two curves.

GENERALIZED FÁRY THEOREM. Let $\gamma, \gamma_0 : \mathbb{S}^1 \rightarrow \mathbb{R}^2$ be two smooth closed curves. Then $|\text{length}(\gamma) - \text{length}(\gamma_0)| \leq [\text{curv}(\gamma) + \text{curv}(\gamma_0) - 2\pi] F(\gamma, \gamma_0)$.

To see that Fáry’s Theorem is indeed a special case, let the image of γ be contained in the unit disk and let the image of γ_0 be a tiny circle centered at the origin, as in Figure VIII.10. Since γ_0 is a circle, its total curvature is 2π . Furthermore, we can make it arbitrarily small so its length approaches zero. While for some curves γ , the Fréchet distance to γ_0 exceeds one, it approaches the maximum distance from the origin, which is at most one. Substituting 0 for $\text{length}(\gamma_0)$, 2π for $\text{curv}(\gamma_0)$, and 1 for $F(\gamma, \gamma_0)$ in the Generalized Fáry Theorem gives the original Fáry Theorem.

Length and total curvature in terms of persistence. A first step toward proving the Generalized Fáry Theorem is a re-interpretation of the length and the total curvature. Fix a direction $u \in \mathbb{S}^1$ and consider $f_u = g_u \circ \gamma$, the height function of the first curve. Almost all level sets, $f_u^{-1}(z)$, consist of an even number of points, decomposing γ into the same number of arcs, half of which belong to the sublevel set, $f_u^{-1}(-\infty, z]$. The number of arcs in the sublevel set is equal to the number of

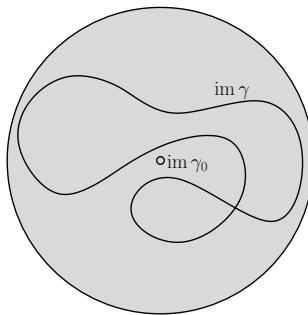


Figure VIII.10: Two curves inside the unit disk. The Fréchet distance between the tiny circle and the other curve approaches a constant at most one as the circle shrinks toward the origin.

components that are born at or before z and are still alive at z . To be precise, this is true as long as z does not exceed the height of the global maximum of f_u . To make it true for all height values, we declare that the component born at the global minimum dies at the global maximum; see Figure VIII.11. This is, incidentally,

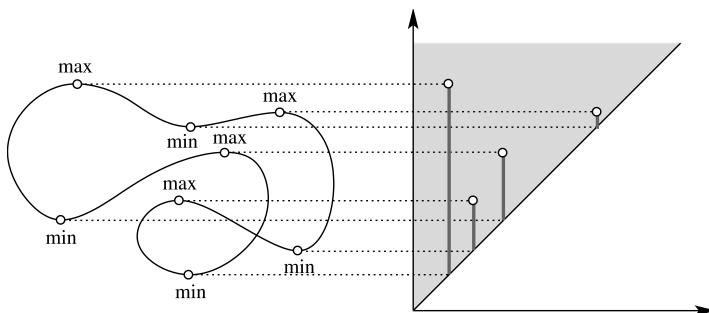


Figure VIII.11: The zeroth persistence diagram of the height function on the curve. We simplify the situation by pairing the global minimum with the global maximum so that all the pair information is contained in this one diagram.

what we would get with extended persistence as described in the previous chapter. Drawing the vertical lines from the off-diagonal points in the persistence diagram down to the diagonal gives a set of line segments with total length equal to the total persistence of $\text{Dgm}_0(f_u)$ as defined in the preceding section. We simplify notation by listing the function rather than its persistence diagram as argument, writing $\Phi(f_u) = \sum \text{pers}(a)$, where the sum is over all points a in $\text{Dgm}_0(f_u)$. By what we said above, the number of line segments that intersect the horizontal line at height z is equal to half the number of points in $f_u^{-1}(z)$. Integrating the number of intersections between γ and lines with normal direction u thus gives twice the total

persistence:

$$\int_{z \in \mathbb{R}} f_u^{-1}(z) = 2\Phi(f_u).$$

The relationship between total curvature and the persistence diagram is even more straightforward. Assuming f_u is Morse, we have a finite number of critical points. This number is even, with equally many minima and maxima paired up to give half the number of off-diagonal points in the persistence diagram. We get similar relationships for the height function of the second curve.

Bounding the difference and integrating. To relate the quantities for the two curves, we write $\varepsilon = F(\gamma, \gamma_0)$ for the Fréchet distance and assume that γ and γ_0 are parametrized such that $\|\gamma(s) - \gamma_0(s)\| \leq \varepsilon$, for all s . It follows that $|f_u(s) - f_{0,u}(s)| \leq \varepsilon$, for all s . The Stability Theorem for Tame Functions then implies that there is a bijection between the points of $\text{Dgm}_0(f_u)$ and of $\text{Dgm}_0(f_{0,u})$ such that corresponding points have L_∞ -distance at most ε . It follows that the difference in persistence between two corresponding points is at most 2ε . If both are off-diagonal points, then we have four critical points (two of f_u and two of $f_{0,u}$) that we can hold responsible for the difference. However, if an off-diagonal point is matched with a point on the diagonal, then we have only two critical points to take responsibility for the 2ε difference. This is indeed the worse of the two possibilities, but we can guarantee that at least two off-diagonal points can be matched within L_∞ -distance ε , namely the two points formed by the global min-max pairs. This is because these critical points correspond to points at infinity in the ordinary persistence diagrams, and being at infinity, they cannot be matched to points on the diagonal. In summary, the difference in total persistence between f_u and $f_{0,u}$ is at most ε times the number of critical points of f_u and $f_{0,u}$ minus two. We are now ready to integrate over all directions $u \in \mathbb{S}^1$ to get the final result. Specifically,

$$\begin{aligned} |\text{length}(\gamma) - \text{length}(\gamma_0)| &\leq \frac{1}{2} \int_{u \in \mathbb{S}^1} |\Phi(f_u) - \Phi(f_{0,u})| du \\ &\leq \frac{\varepsilon}{2} \int_{u \in \mathbb{S}^1} [\#\text{crit}(f) + \#\text{crit}(f_0) - 2] du \\ &= \varepsilon[\text{curv}(\gamma) + \text{curv}(\gamma_0) - 2\pi], \end{aligned}$$

using first the Cauchy-Crofton Formula, second the re-interpretations in terms of persistence, third the inequality implied by the Stability Theorem for Functions, and fourth the Total Curvature Formula. This completes the proof of the Generalized Fáry Theorem.

Bibliographic notes. The inequality that connects the length with the total curvature of a closed curve is due to Fáry [67]. The generalization that compares the lengths of curves that are close in the Fréchet distance is more recent [33]. Both results have generalizations to curves in dimensions beyond two. The integral geometry interpretations of length and total curvature can be found in Santaló [129].

VIII.4 Bipartite Graph Matching

In this section, we consider algorithms for the bottleneck and Wasserstein distances between persistence diagrams. Both problems reduce to constructing optimal matchings in bipartite graphs.

Distance from matching. We begin by reducing the computation of distance to constructing a matching. Let X and Y be two persistence diagrams. We assume both consist of finitely many points above the diagonal and infinitely many points on the diagonal. Letting X_0 be the finite multiset of off-diagonal points in X and X'_0 the orthogonal projection of X_0 onto the diagonal, we construct a complete bipartite graph $G = (U \dot{\cup} V, E)$ with $U = X_0 \dot{\cup} Y'_0$, $V = Y_0 \dot{\cup} X'_0$, and $E = U \times V$. For each $q > 0$, we introduce the cost function $c = c^q : E \rightarrow \mathbb{R}$ defined by mapping the edge $uv \in E$ to the q -th power of the L_∞ -distance between the points:

$$c(uv) = \begin{cases} \|u - v\|_\infty^q & \text{if } u \in X_0 \text{ or } v \in Y_0; \\ 0 & \text{if } u \in Y'_0 \text{ and } v \in X'_0. \end{cases}$$

By construction, the minimum cost edge connecting an off-diagonal point u to a point on the diagonal is the edge uu' , where u' is the orthogonal projection of u . For $q = 1$, the cost of this edge is half the persistence of u .

A *matching* of G is a subset of vertex disjoint edges, $M \subseteq E$. It is *maximum* if there is no matching with more edges and *perfect* if every vertex is the endpoint of an edge in M . Since G is complete with equally many vertices on the two sides, every maximum matching is also a perfect matching. We will also consider matchings for graphs $G(\varepsilon) = (U \dot{\cup} V, E_\varepsilon)$ obtained from G by removing all edges $uv \in E$ with cost $c(uv) > \varepsilon$. Of course, every perfect matching of $G(\varepsilon)$ is a maximum matching but not necessarily the other way around. A *minimum cost matching* is a maximum matching that minimizes the sum of costs of the edges in the matching. We refer to this sum as the *total cost* of the matching. It is not difficult to prove the following relation between distance and matching.

REDUCTION LEMMA. Let X and Y be two persistence diagrams and let $G = (U \dot{\cup} V, E)$ be the corresponding complete bipartite graph. Then

- (i) the bottleneck distance between X and Y is the smallest $\varepsilon \geq 0$ such that the subgraph $G(\varepsilon)$ of G with cost function $c = c^1$ has a perfect matching;
- (ii) the q -th Wasserstein distance between X and Y is the q -th root of the total cost of the minimum cost matching of G with cost function $c = c^q$.

We are therefore interested in recognizing bipartite graphs that have perfect matchings and in constructing minimum cost matchings.

Augmenting paths. We begin by considering the algorithmic problem of constructing a maximum matching of the bipartite graph $G(\varepsilon) = (U \dot{\cup} V, E_\varepsilon)$. The

algorithm is iterative, improving the matching in each round, until no further improvement is possible. Let M_i be the matching after i iterations. The crucial concept is a path that alternates between edges in and out of M_i . To explain this, we introduce a directed graph D_i that depends on $G(\varepsilon)$ and M_i . For the most part, it is the same as $G(\varepsilon)$ except that each edge is drawn with a direction, namely from V to U if the edge belongs to M_i , and from U to V if the edge does not belong to M_i . In addition to the vertices in $G(\varepsilon)$, the directed graph contains two new vertices, the source s with an edge from s to every unmatched vertex $u \in U$, and the target t with an edge from every unmatched vertex $v \in V$ to t ; see Figure VIII.12. An *augmenting path* is a directed path from s to t that visits every vertex

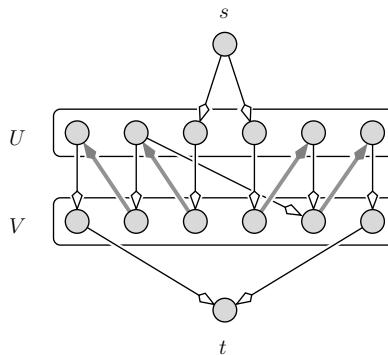


Figure VIII.12: A bipartite graph with six plus six vertices and a matching with four edges giving rise to a directed graph with three paths from s to t .

at most once. By construction, an augmenting path consists of $2k + 1$ edges, one from s to U , an interleaved sequence of k edges not in M_i and $k - 1$ edges in M_i , and finally an edge from V to t . Clearly, if we have an augmenting path, we can improve the matching by substituting the k edges not in M_i for the $k - 1$ edges in M_i . When we make this improvement, we say we *augment* the matching using the path. To get an algorithm, we also need the existence of an augmenting path unless M_i is maximum. To construct such a path, draw the edges of an assumed maximum matching from U to V and those of M_i from V to U . Each vertex is incident to at most two edges, one incoming and the other outgoing, so we can partition the edges into maximal, vertex disjoint paths and closed curves that interleave edges from the two matchings. A path in this partition extends to an augmenting path from s to t iff it contains one more edge from the maximum matching than from M_i . Since M_i is smaller, there is at least one such path. We use this fact to give an algorithm for constructing a maximum matching of $G(\varepsilon)$.

```

 $M_0 = \emptyset; i = 0;$ 
while there exists an augmenting path in  $D_i$  do
    augment  $M_i$  using this path to get  $M_{i+1}$ ;
     $i = i + 1$ 
endwhile.

```

Each iteration increases the size of the matching by one. The number of edges in the maximum matching is at most $n = \text{card } U = \text{card } V$, which implies that the algorithm terminates after at most n iterations. We can use Depth-first Search or Breadth-first Search to find an alternating path in time proportional to the number of edges, $m_\varepsilon = \text{card } E_\varepsilon$. In either case, we have an algorithm that runs in time at most proportional to $m_\varepsilon n \leq n^3$.

Shortest augmenting paths. The running time of the algorithm can be improved if we use multiple augmenting paths at a time. Specifically, we use a maximal set of edge disjoint, shortest, augmenting paths. To find them, we use Breadth-first Search to label all vertices by their distance from the source and Depth-first Search to construct a maximal set of paths in the thus labeled directed graph. Since Depth-first Search has been explained in detail in Section II.2, we focus on the first step.

```

 $S_0 = \{s\}$ ; label  $s$  with 0;  $j = 0$ ;
while  $S_j \neq \emptyset$  do
    forall vertices  $x \in S_j$  do
        forall unlabeled successors  $y$  of  $x$  do
            label  $y$  with  $j + 1$  and add  $y$  to  $S_{j+1}$ 
        endfor
    endfor;  $j = j + 1$ 
endwhile.

```

Assuming suitable data structures, we can iterate through the vertices in the sets S_j and their successors in constant time per vertex. Using repeated Depth-first Search in the labeled graph D_i , we construct a maximal set of edge disjoint paths from s to t . If we remove edges and vertices as they become useless, we get an algorithm that computes the paths in time proportional to m_ε . For example, if we start with the directed graph in Figure VIII.12, we get either two paths of length seven, one on the left and the other on the right, or just one path of the same length, as shown in Figure VIII.13. Finally, we augment the matching using all paths in the maximal set.

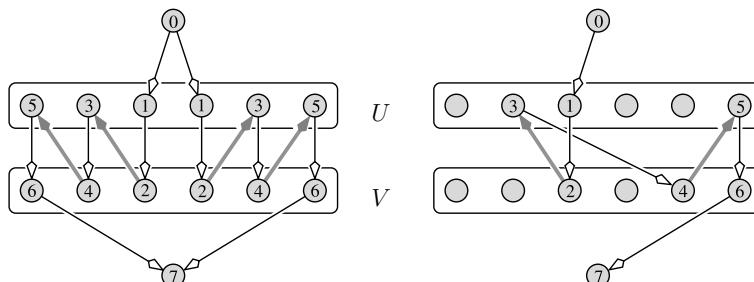


Figure VIII.13: The two maximal sets of edge disjoint, shortest, augmenting paths in the directed graph of Figure VIII.12.

Analysis. We now show that the new strategy leads to a substantially smaller number of iterations. In a nutshell, the reason is that there cannot be many augmenting paths that are all long. Playing off length against number, we get a bound of some constant multiplied with the square root of the number of vertices.

ITERATION BOUND. Starting with the empty matching and augmenting the matching of $G(\varepsilon)$ using a maximal set of edge disjoint, shortest, augmenting paths each time, we reach a maximum matching in fewer than $2\sqrt{2n}$ iterations.

PROOF. We first show that the length of the shortest path from s to t increases from one iteration to the next. Let $\ell_i(x)$ be the length of the shortest path from s to the vertex x in D_i ; it is the label assigned to x by Breadth-first Search. We prove that $\ell_{i+1}(t)$ is strictly larger than $\ell_i(t)$, assuming both are defined. Consider a shortest path π from s to t in D_{i+1} . It is also a path in D_i iff none of its edges belongs to the paths selected in the i -th round. If π is a path in D_i , then it cannot be shortest; otherwise, it would have been added to the maximal set. On the other hand, if π is not in D_i , then it has at least one edge xy that is reversed in D_i . Since yx belongs to a shortest path in D_i , we have $\ell_i(y) = \ell_i(x) - 1$. For an edge xy of π that is not reversed in D_i , we have $\ell_i(y) \leq \ell_i(x) + 1$ by definition of ℓ_i . As we walk along the path, ℓ_{i+1} grows by one at each step while ℓ_i grows by at most one and at least once it shrinks. Hence $\ell_i(t) < \ell_{i+1}(t)$, as required.

For the second part of the proof, we note that two edge disjoint paths from s to t share no vertices other than the source and the target. This is because each vertex of U has only one incoming edge and each vertex of V has only one outgoing edge. Let \bar{m}_i be the size deficit of M_i ; that is, the number of edges is short of being a maximum matching. Since M_i can be improved by this much, there are at least \bar{m}_i augmenting paths from s to t in D_i . Using the construction of augmenting paths given earlier in this section, we find \bar{m}_i augmenting paths that share no vertices other than s and t . By the pigeonhole principle, the shortest of these paths contains at most a fraction of $1/\bar{m}_i$ of the vertices of $G(\varepsilon)$. Equivalently, $\ell_i(t) \leq 2n/\bar{m}_i + 1$. Since the distance of t from s begins at three and grows with increasing i , this implies $i \leq 2n/\bar{m}_i - 2$. To increase M_i by another \bar{m}_i edges takes at most \bar{m}_i additional iterations. The total number of iterations is therefore bounded from above by $2n/\bar{m}_i - 2 + \bar{m}_i$. Setting \bar{m}_i to the smallest integer no smaller than $\sqrt{2n}$ implies the claimed bound. \square

Recall that each iteration takes time at most proportional to the number of edges. The bound on the number of iterations thus implies that the algorithm runs in time at most proportional to $m_\varepsilon \sqrt{n} \leq n^{5/2}$.

Minimum cost matching. To compute the smallest ε for which $G(\varepsilon)$ has a perfect matching, we do binary search in the list of edges sorted by cost, constructing a maximum matching at every step. Similarly, constructing a minimum cost matching of G is done by iterating the maximum matching algorithm, but the iteration is different. There are two easy structural insights that show the way.

1. If the subgraph $G(0)$ consisting of the cost zero edges in G has a perfect matching, then this is a minimum cost matching. Indeed, its total cost is zero, which is as small as it gets.
2. Subtracting the same amount from the cost of all edges incident to a vertex in G affects all perfect matchings the same way. In particular, a perfect matching minimizes the total cost before the subtractions iff it does so after the subtractions.

To compute a minimum cost matching of G , we begin with all zero cost edges and construct a maximum matching of $G(0)$. If the matching is perfect, we are done. Otherwise, we change the costs of the edges in G while preserving the ordering of the perfect matchings by total cost. To describe how this is done, we introduce *deduction maps* $d_i : U \cup V \rightarrow \mathbb{R}$. Starting with the zero map, $d_0(x) = 0$ for all vertices x , the algorithm will change the map and will this way modify the costs. Writing $c(xy)$ for the original cost of the edge xy in G , the *modified cost* after i iterations is

$$c_i(xy) = c(xy) - d_i(x) - d_i(y).$$

It is important for the efficiency but also the correctness of the algorithm that all modified costs always be non-negative. This will be an invariant of the algorithm. Letting G_i be the graph G with costs modified using d_i , the algorithm iterates the construction of a maximum matching of $G_i(0)$, the graph G_i with edges of positive modified cost removed. Increasing the maximum matching by one edge each time, we get a perfect matching after n iterations. By construction, all edges in this matching have zero modified cost.

Minimum cost paths. We now show how to change the deduction map so that the maximum matching increases. Let M_i be a maximum matching of $G_i(0)$ and let $D_i(0)$ be the directed graph defined by $G_i(0)$ and M_i . Because M_i is maximum, $D_i(0)$ has no directed path from s to t . Let D_i be the directed graph defined by G_i and the same matching M_i and note that it contains $D_i(0)$ as a subgraph. Assuming M_i is not perfect, it is not maximum for G_i , which implies that D_i has directed paths from s to t . Each such path is an augmenting path, and we define its *total cost* as the sum of modified costs of its edges. By definition, the modified cost of the first edge, from s to U , is zero, and so is the modified cost of the last edge, from V to t . Let π be the augmenting path in D_i that minimizes the total cost. It can be computed by an algorithm similar to Breadth-first Search. Indeed, the only difference is that it visits the vertices in a particular ordering that depends on the modified costs of the edges. At every moment during the construction, we have a set of visited vertices forming a tree rooted at s and we have a set of unvisited vertices. For each unvisited vertex, y , we consider the minimum cost path that starts at s , goes to a vertex x using edges in the tree, and ends with the edge from x to y . The next vertex visited by the algorithm is the unvisited vertex y that minimizes this cost, and we add y together with the last edge of its path to the tree. This is known as Dijkstra's Single Source Shortest Path Algorithm, or Dijkstra's

Algorithm for short. We compute the minimizing vertex y and update the costs of all yet unvisited vertices in time proportional to n . Iterating this step n times, we find the minimum cost path π in time proportional to n^2 .

We augment the matching M_i using π to get M_{i+1} . This increases the matching, but to be sure that we made progress toward computing a minimum cost matching, we have to show that it is possible to change the deduction map so that all edges in M_{i+1} have zero modified costs. To this end, let $\gamma_i(x)$ be the minimum total cost of a path from s to x ; it is the total cost of the path from s to x within the tree computed by Dijkstra's Algorithm. Using these quantities, we update the deduction map to

$$d_{i+1}(x) = \begin{cases} d_i(x) - \gamma_i(x) & \text{if } x \in U; \\ d_i(x) + \gamma_i(x) & \text{if } x \in V. \end{cases}$$

For vertices $u \in U$ and $v \in V$, the new modified cost of the edge connecting u with v is

$$\begin{aligned} c_{i+1}(uv) &= c(uv) - d_{i+1}(u) - d_{i+1}(v) \\ &= c(uv) - d_i(u) - d_i(v) + \gamma_i(u) - \gamma_i(v). \end{aligned}$$

In words, it is the old modified cost plus $\gamma_i(u) - \gamma_i(v)$, no matter whether the edge goes from u to v or from v to u in D_i . If $\gamma_i(u) \geq \gamma_i(v)$, we use induction to get $c_{i+1}(uv) \geq 0$ from $c_i(uv) \geq 0$. Otherwise, $\gamma_i(v) - \gamma_i(u) \leq c_i(uv)$, because $\gamma_i(v)$ is the minimum total cost of a path from s to v . It follows that all new modified costs are non-negative. But we need more, namely zero new modified cost for all edges of the new matching. There are two kinds of such edges uv , those that belong to M_i and those that belong to the path π . For the first kind, we have $\gamma_i(v) = \gamma_i(u)$ because $c_i(uv) = 0$, and the only way to reach u is along the directed edge from v to u . For the second kind, we have $\gamma_i(v) - \gamma_i(u) = c_i(uv)$ by definition of γ_i . In both cases, we have $c_{i+1}(uv) = 0$, as required.

This completes the proof that the iteration ends with a perfect matching minimizing the total cost. The maximum matching gains one edge per iteration. We thus have n iterations each taking time proportional to n^2 . Our algorithm thus constructs a minimum cost matching in time at most proportional to n^3 .

Bibliographic notes. Computing a maximum matching of a bipartite graph is a classic optimization problem discussed in operations research texts [7]. As explained in [140], it is a special case of the more general maximum flow problem in networks. Indeed, Dinic's maximum flow algorithm for so-called unit networks [50] specializes to the $n^{5/2}$ time algorithm for maximum matching independently discovered by Hopcroft and Karp [85] and explained in this section. The Minimum Cost Matching Algorithm is a variant of what is known as the Hungarian method [96]. Following [92], we describe a version that uses Dijkstra's Algorithm for finding shortest paths in a weighted graph as a subroutine [49]. Using the geometry of the persistence diagrams, the Maximum Matching Algorithm can be improved to run in time at most proportional to $n^{3/2} \log_2 n$ [65], and the Minimum Cost Matching Algorithm can be improved to run in time at most proportional to $n^{2+\varepsilon}$ [5].

Exercises

The credit assignment reflects a subjective assessment of difficulty. A typical question can be answered using knowledge of the material combined with some thought and analysis.

1. **Examples of switches** (two credits). Give examples for the types of switches analogous to the ones shown in Figure VIII.4 but one dimension up in each of the three types.
2. **Matrix maintenance** (two credits). Formulate an algorithm that maintains the reduced boundary matrix under transpositions for $\partial = RV$; that is, it maintains the matrix V instead of its inverse, U .
3. **Sparse matrix representation** (two credits). Give a sparse matrix representation that allows an implementation of the maintenance algorithm running in time proportional to the number of ones in the changed columns and rows of $R = \partial U$.
4. **Measuring vineyards** (two credits). Let $f, g : K \rightarrow \mathbb{R}$ be two monotonic functions on a simplicial complex and let $f_t = (1-t)f + tg$ for $t \in [0, 1]$ be the straight-line homotopy between them. Each vine of the homotopy is a map $x : [a, b] \rightarrow \bar{\mathbb{R}}^2$ with $0 \leq a < b \leq 1$. Let

$$\mu(x) = \int_{s=a}^b \|x(s) - x(a)\| ds$$

and define a measure by summing the integrals over all vines in the p -th vineyard, $\mu_p(f, g) = \sum_x \mu(x)$. Give examples that show that μ_p and the first Wasserstein distance are incomparable, that is, there are monotonic functions f, g, f_0, g_0 such that $\mu_p(f, g) < W_1(\text{Dgm}_p(f), \text{Dgm}_p(g))$ and $\mu_p(f_0, g_0) > W_1(\text{Dgm}_p(f_0), \text{Dgm}_p(g_0))$.

5. **Cauchy-Crofton** (two credits). Generalize the Cauchy-Crofton formula for curves in the plane given in Section V.3 to
 - (i) curves in 3-dimensional Euclidean space;
 - (ii) surfaces in 3-dimensional Euclidean space.

6. **Mean and Gaussian curvatures** (three credits). Use the structure of the proof of the Generalized Fáry Theorem to show the following relationship between the total mean curvature and the total absolute Gaussian curvature of two homeomorphic closed surfaces embedded in \mathbb{R}^3 :

$$|\text{mean}(S) - \text{mean}(S_0)| \leq [\text{gauss}(S) + \text{gauss}(S_0) - 4\pi(1+g)]F(\bar{S}, \bar{S}_0),$$

where g is the common genus of S and S_0 , \bar{S} and \bar{S}_0 are the solid bodies bounded by the two surfaces, and $F(\bar{S}, \bar{S}_0)$ is the Fréchet distance between them.

7. **Breadth-first Search** (one credit). Reformulate the Breadth-first Search algorithm for labeling the vertices of D_i using a single queue to represent all sets of vertices S_j in one data structure. As suggested by the name, this is a data structure that supports adding an element at the end and removing it from the front, both in constant time.
8. **Incremental matching** (three credits). Recall that the maximum matching of a bipartite graph with n vertices can be constructed in time at most proportional to $n^{5/2}$. Running this algorithm within a binary search routine, we find the perfect matching of a complete bipartite graph that minimizes the largest cost of any of its edges in time at most proportional to $n^{5/2} \log_2 n$. Show that the two algorithms can be integrated to avoid the $\log_2 n$ overhead, constructing the perfect matching in time at most proportional to $n^{5/2}$.

Chapter IX

Applications

The primary application of the mathematical and computational tools introduced in the previous chapters is in data analysis, an activity that reaches into every discipline in science and engineering. The data may comprise the readings of an array of sensors, the pixels of an image, the accumulation of observations, or what have you. Invariably, there is noise in the data, which may be systematic or random. It may also reflect genuine properties of the measured phenomenon but at a scale that is outside the window of interest. The traditional approach to noise is to ‘smooth’ or ‘regularize’ the data, which invariably means we change the data. This is in contrast to the approach we advocate here, namely to measure the noise and not change the data. What is new is the measurement and the additional level of rationality and consistency it affords us. The four case studies selected to illustrate the possibilities that this paradigm affords us all start with biological data.

IX.1 Measures for Gene Expression Data

Our first application deals with 1-dimensional real-valued functions, the simplest kind of objects about which persistent homology can make meaningful statements. Such functions arise in the development of somites in vertebrates.

Background. Vertebrates are characterized by a spinal column consisting of a sequence of vertebrae that provide a periodic segmentation of their body along the axis. Mice are one example, with a spinal column of about sixty-five vertebrae. The numbers are larger for snakes, whose columns might be segmented into a few hundred vertebrae. This structure arises in the development of the embryo, when the vertebral precursors, the *somites*, are formed rhythmically from the presomitic mesoderm. This process is associated with a molecular oscillator that drives gene expression with a period corresponding to that of somite formation. We refer to this oscillator as the *segmentation clock*. The desire to fully understand this clock is the motivation for the work described in this section.

An early indication of the molecular underpinnings of somite development was the visual exposure of a cyclically expressed gene called lunatic fringe. Adding a fluorescent marker, its expression could be observed as a wave initiated in the posterior presomitic mesoderm. Migrating up, the wave narrows as it moves to the anterior, where the somites form.

Technology. The segmentation clock is one of the most reliable organic structures, and it has a built-in counter that terminates its rhythm after some number of periods. Its operation suggests an elaborate mechanism involving more than a few genes. Microarray technology offers a way to pursue the broad question of which genes are involved by testing the entire genome of an organism at once. An array is a 2-dimensional organization of array elements, each measuring the expression of a particular gene. This is done by depositing pieces of DNA that are specific to the RNA product of that particular gene. These pieces bind to copies of particular RNA strands, if they are present in the tissue probe. The binding event is made observable by fluorescence whose intensity quantifies the abundance of the particular strand in the tissue.

The organism of choice for the study of the segmentation clock is the mouse. We start with a microarray primed with the entire mouse genome. Copies of this array are used to measure the expression of all genes several times during a single period. In the mouse, a somite is developed roughly every two hours, and measurements are taken at seventeen time points in that interval. It is important to mention that this description is a simplification of the actual experiment. Tissue probes are taken from seventeen embryos and during five periods. Rather than timing the probes with a stopwatch, the time within a period is estimated from the state of the observed wave of lunatic fringe expression. Instead of quantified time, we thus have ranked time, events subjectively ordered by visual inspection. In the end, we have a series of seventeen measurements for each of about seven and a half thousand genes in the mouse genome. Each measurement is a real number representing the observed intensity at the particular array location, which quantifies the abundance of the corresponding strand of RNA.

Before discussing the mathematical analysis of this data, we draw attention to an inherent limitation that results from folding data from several periods into one. Suppose we have a gene that is rhythmically expressed but with a different period, say, three instead of two hours. The sorting process will shuffle the data collected for this gene, destroying any clear signal if there was one. It is thus reasonable to use this data to decide whether a gene is rhythmically expressed with a period consistent with somite development, but not whether a gene is rhythmically expressed at all, or what the most likely length of the period would be.

Lipschitz functions on the circle. We model the results of the time series of microarray experiments as a set of functions from the circle to the real numbers, $f : \mathbb{S}^1 \rightarrow \mathbb{R}$, one for each gene. The circle represents the two hours of one period, and the function tracks the abundance of the gene product within the period. Change requires energy, namely for the production and degradation of RNA. We use this

as a justification to assume that f is Lipschitz, that is, there is a smallest positive constant, called the *Lipschitz constant* of f and denoted as $\text{Lip}(f)$, such that $|f(s) - f(t)| \leq \text{Lip}(f)\|s - t\|$, where the distance between s and t is measured along the circle. For a differentiable function, this is equivalent to constraining the derivative between $\pm \text{Lip}(f)$. Defining the *total variation* as the integral of the norm of the derivative, we thus get

$$\text{Var}(f) = \int_{s=0}^{2\pi} |f'(s)| ds \leq 2\pi \text{Lip}(f).$$

This inequality will be relevant shortly, when we study the stability of different ways to measure functions on the circle. To prepare this study, we consider the persistence diagram of f . It expresses the history of births and deaths in the sequence of sublevel sets, $f^{-1}(-\infty, a]$. Assuming f is Morse, we have the birth of a component at every minimum and the death of a component at every maximum, except at the last, global maximum at which we have the birth of a 1-dimensional class. This class never dies. Similarly, the 0-dimensional class born at the global minimum never dies. All other classes are born and die at finite values. Using the notation from Chapter VI, we write c_0 for the number of minima and c_1 for the number of maxima of f . Clearly, $c_0 = c_1$. By what we said above, the only persistence diagram that contains interesting information is the zeroth, $\text{Dgm}_0(f)$, containing one point at infinity and $n = c_0 - 1 = c_1 - 1$ points in its finite portion, \mathbb{R}^2 . Each finite point corresponds to a minimum paired with a maximum, and we write $x_i = (b_i, d_i)$, where b_i and d_i are the values of f at the minimum and the maximum. Let b_0 and b_{n+1} be the values of f at the global minimum and the global maximum, remembering that (b_0, ∞) is the point at infinity in $\text{Dgm}_0(f)$ and (b_{n+1}, ∞) is the only point in $\text{Dgm}_1(f)$. Consistent with the notation in the preceding chapter, we write

$$\Phi(f) = \sum_{i=1}^n (d_i - b_i)$$

for the total persistence of f . Note that this is the same as half the total variation minus the amplitude; that is, $\Phi(f) = \frac{1}{2}\text{Var}(f) - (b_{n+1} - b_0)$. Indeed, $\Phi(f) + (b_{n+1} - b_0)$ is the sum of the values of f at the $n + 1$ maxima minus the sum of the values at the $n + 1$ minima. Decomposing f into increasing and decreasing portions, we can write $\text{Var}(f)$ as the sum of two integrals, each equal to the same difference of sums. Hence, $\text{Var}(f) = 2\Phi(f) + 2(b_{n+1} - b_0)$, as claimed.

Simplification. Before we introduce a measure for how close a function $f : \mathbb{S}^1 \rightarrow \mathbb{R}$ is to being periodic, in our assessment, we need to understand how we can simplify. Call another continuous function $f_\varepsilon : \mathbb{S}^1 \rightarrow \mathbb{R}$ an ε -simplification of f if

- (i) $|f(s) - f_\varepsilon(s)| \leq \varepsilon$ for all $s \in \mathbb{S}^1$;
- (ii) an off-diagonal point belongs to $\text{Dgm}_p(f_\varepsilon)$ iff it belongs to $\text{Dgm}_p(f)$ and its vertical distance from the diagonal exceeds ε .

Condition (ii) says the persistence diagrams of the two functions are the same except that the diagrams of f_ε have no points of persistence ε or less. We prove the existence of ε -simplifications by explicit construction of a function f_ε . It is convenient to

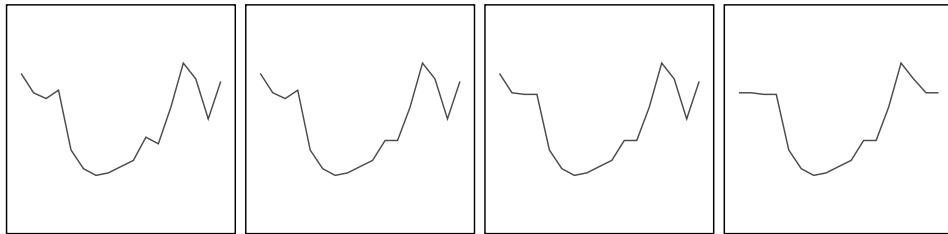


Figure IX.1: Simplifications of the expression profile of the gene Axin2. From left to right: the original function and the simplified functions obtained by canceling one, two, and all three minimum-maximum pairs. Notice that the last cancellation affects the curve at both ends, because the domain of the function is the circle.

assume that f is PL Morse. If f has only one minimum and one maximum, then its persistence diagrams have no finite points and f is its own ε -simplification for every real number $\varepsilon \geq 0$. So assume f has at least two minima and two maxima. Let u and v be a pair with minimum persistence and let $x = (f(u), f(v))$ be the corresponding point in the zeroth diagram. By assumption of minimality, the function value increases monotonically to $c = f(v)$ on both sides of u . Let $t \neq v$ be the point with $f(t) = c$ reached from u going in the direction away from v . If $f(v) - f(u) < \varepsilon$, we change f by setting $f(s) = c$ for all points s on the arc from t to v that contains u . The values on the complementary arc are preserved. We can make the new function PL Morse by giving a subtle slope to the flat interval between v and t , slightly extending it beyond t to pick up a small amount of height. The persistence diagrams of the new f are the same as before, except that the point x has disappeared. We get f_ε by repeating this step for all minimum-maximum pairs with persistence ε or less. This construction is illustrated in Figure IX.1, which shows the function for the gene Axin2 along with three simplifications.

Measures. The sine function, which maps points of \mathbb{S}^1 to their second Cartesian coordinates in \mathbb{R}^2 , is the prototypical periodic function. It has a single minimum and a single maximum and varies smoothly between the two. Allowing for more general patterns to increase and decrease, we retain the property of having only two critical points as the characteristic ideal of a periodic function. To quantify periodicity more generally, we assign zero to a function with $c_0 + c_1 = 2$ and a positive number to every other function. Again, we find it convenient to restrict the discussion to PL Morse functions. Specifically, we set $\mu_0(f) = \frac{1}{2}(c_0 + c_1) - 1$, and for every positive integer q , we define the *degree- q periodicity measure* by integrating the degree- $(q - 1)$ measure over the ε -simplifications of f :

$$\mu_q(f) = \int_{\varepsilon \geq 0} \mu_{q-1}(f_\varepsilon) d\varepsilon.$$

Note that $\mu_1(f)$ is proportional to the average number of critical points of the ε -simplifications. To see that $\mu_q(f)$ is well defined, we show that the measures do not depend on which ε -simplifications we use. We do this by proving that μ_q is equal to the *degree- q total persistence* of f defined as $\Phi^q(f) = \sum_{i=1}^n (d_i - b_i)^q$.

PERIODICITY MEASURE LEMMA. Let $f : \mathbb{S}^1 \rightarrow \mathbb{R}$ be a PL Morse function. Then $\mu_q(f) = \Phi^q(f)$ for all non-negative integers q .

PROOF. We use induction over q to prove that the contribution of the point $x_i = (b_i, d_i)$ in $\text{Dgm}_0(f)$ to the degree- q periodicity measure is $(d_i - b_i)^q$. For $q = 0$, this point contributes one to $\mu_0(f)$ as well as to $\Phi^0(f)$. This establishes the base case. Let $q \geq 1$. By definition of ε -simplification, the point x_i belongs to the zeroth diagram of f_ε for all $0 \leq \varepsilon < d_i - b_i$ but not for any larger values of ε . The contribution of x_i to the degree- q periodicity measure is therefore $d_i - b_i$ times its contribution to the degree- $(q-1)$ measure, which, by inductive assumption, is $(d_i - b_i)^{q-1}$. Summing over all points x_i , for $1 \leq i \leq n$, gives $\Phi^q(f)$. There are no other finite points in the diagrams of the f_ε , which implies the claim. \square

The Periodicity Measure Lemma implies an algorithm for computing $\mu_q(f)$, namely constructing the zeroth persistence diagram and summing the q -th powers of the vertical distances of its finite points from the diagonal. It also provides a definition of μ_q for real values q that are not integers.

Instability for small degree. Whether or not the periodicity measure is stable depends on the choice of q . Clearly, μ_0 is not stable because arbitrarily small perturbations can change the measure by an arbitrary amount. Perhaps less obviously, μ_1 is also not stable. Perhaps less obviously, μ_1 is also not stable. To see this,

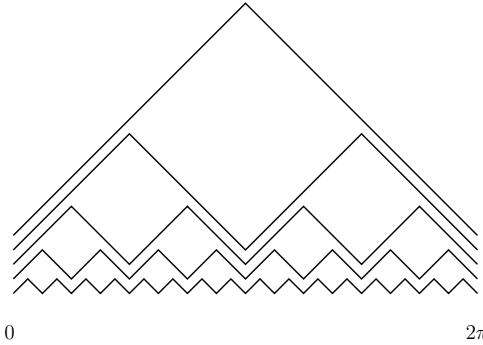


Figure IX.2: The graphs of the functions g_k , for $k = 0, 1, 2, 3, 4$, with vertical off-set for clarity.

we construct a series of Lipschitz functions, $g_k : \mathbb{S}^1 \rightarrow \mathbb{R}$, that approach the zero function while their total persistence approaches π . Replacing each point of \mathbb{S}^1 by its angle, $\varphi \in [0, 2\pi]$, we set $g_0(\varphi) = \min\{\varphi, 2\pi - \varphi\}$ and define $g_k(\varphi) = \frac{1}{2}g_{k-1}(2\varphi)$ for all positive integers k ; see Figure IX.2. The maximum difference between g_k and

the zero function is $\|g_k\|_\infty = \max_{0 \leq \varphi < 2\pi} g(\varphi) = \pi/2^k$, which goes to zero as k goes to infinity. On the other hand, every function g_k has slope ± 1 almost everywhere. The total variation is therefore $\text{Var}(g_k) = 2\pi$. We divide by two and subtract the amplitude to get the total persistence as $\Phi^1(f) = \pi - \pi/2^k$, which goes to π as k goes to infinity.

Stability for degree at least two. There is a qualitative difference between the periodicity measures when q passes from one to two. In particular, μ_q is stable for every constant $q \geq 2$.

STABILITY THEOREM FOR TOTAL PERSISTENCE. Let $f, g : \mathbb{S}^1 \rightarrow \mathbb{R}$ be Lipschitz functions with Lipschitz constant one and let $q \geq 2$. Then

$$|\Phi^q(f) - \Phi^q(g)| \leq 4q\pi^{q-1} \cdot \|f - g\|_\infty.$$

PROOF. We begin by noting that $y^q - x^q = \int_x^y qt^{q-1} dt \leq q|y - x| \max\{x, y\}^{q-1}$ for all $x, y \geq 0$ and $q \geq 1$. We use the Stability Theorem for Tame Functions in Chapter VIII to index the persistences of the finite points in the zeroth diagrams of f and g such that

$$\begin{aligned} \Phi^1(f) &= \phi_1 + \phi_2 + \dots + \phi_m, \\ \Phi^1(g) &= \gamma_1 + \gamma_2 + \dots + \gamma_m, \end{aligned}$$

and $|\phi_i - \gamma_i| \leq 2\varepsilon$ for all i , where $\varepsilon = \|f - g\|_\infty$, possibly after adding zeros. Both sums are bounded from above by half the total variation, which implies $\Phi^1(f) + \Phi^1(g) \leq 2\pi$. We also note that $\phi_i \leq \pi$ and $\gamma_i \leq \pi$ for $1 \leq i \leq n$. Writing $\Delta = \Phi^q(f) - \Phi^q(g)$, we therefore get

$$\begin{aligned} |\Delta| &\leq \sum_{i=1}^m |\phi_i^q - \gamma_i^q| \\ &\leq \sum_{i=1}^m q|\phi_i - \gamma_i| \max\{\phi_i, \gamma_i\}^{q-1} \\ &\leq q(2\varepsilon)\pi^{q-2} \sum_{i=1}^m \max\{\phi_i, \gamma_i\}. \end{aligned}$$

The sum in the last expression is bounded from above by $\sum_{i=1}^m (\phi_i + \gamma_i) \leq 2\pi$. The claimed inequality follows. \square

For constant $q \geq 2$, the right-hand side of the inequality in the theorem is at most some constant times the L_∞ -difference between the two functions. It follows that the difference between the degree- q total persistences goes to zero as the difference between the functions goes to zero. The above theorem is thus a statement of stability for total persistence and therefore for the periodicity measure.

Notes. The background for the material in this section is provided by the biological work on somite development in Pourquié's group; see e.g. [118, 123]. The microarray time series data of the mouse genome forms the motivation for our work. It had originally been analyzed using a variant of Fourier analysis [47]. Because of limitations in the discerned patterns, the same data was later re-analyzed using four custom-made mathematical methods, all designed to recognize rhythmic gene expression of the kind exhibited by a small number (fewer than 30) of genes verified to participate in somite development. One of these methods was the periodicity measure described in this section. Assessing all seven and a half thousand expression profiles, each method generated a list of the genes, ordered from most to least compatible with the rhythm of somite development. These lists were then compared on the basis of their ranking of the verified genes. The results of the comparison can be found in [46], including the discussion of a small number of newly identified genes.

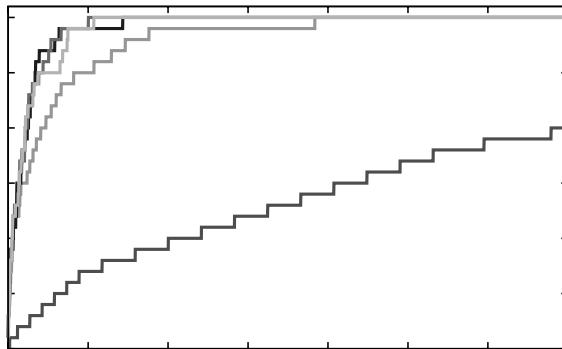


Figure IX.3: The step functions characterizing the distribution of the verified genes in the ordered lists generated with the periodicity measures μ_q , for $q = 0, 1, 2, 3, 4$. For better visibility, we truncate the lists after the first three and a half thousand genes. Moving northwest, toward the upper left corner of the rectangle, we first cross the graph of f_0 , then that of f_1 , and finally the graphs of f_2, f_3, f_4 in an order that depends on the exact route we choose.

The particular periodicity measure used in the re-analysis was μ_2 , which we preferred over the other choices because $q = 2$ is the smallest power for which we know that the measure is stable. There is indeed evidence that stability is an important property for the task at hand. This is illustrated in the following comparison of the measures for $q = 0, 1, 2, 3, 4$. For each q , we generate an ordered list of the genes, as before, and we construct a step function, $f_q : [0, 1] \rightarrow [0, 1]$, that counts the verified genes in every initial segment of the list. In other words, $f_q(x)$ equals the percentage of the total number of verified genes that lie within the initial x percent of the list. If the verified genes are distributed evenly among the others, then we get a step function whose graph is close to the diagonal. On the other hand, if the verified genes are all listed first, then the function shoots up to one and stays there until the end. In general, one measure performs better than another if the first function majorizes the second. As shown in Figure IX.3, there

is indeed a marked difference between the unstable measures, μ_0 and μ_1 , and the stable measures, μ_2 , μ_3 , and μ_4 . In summary, the graph of f_0 is slightly above the diagonal, indicating that μ_0 performs only marginally better than giving a random ordering. The visibly most striking improvement is from μ_0 to μ_1 . However, as we get closer to the ideal step function, improvements are more difficult to come by, so the improvement from μ_1 to μ_2 is also significant. Thereafter, the graphs for $q = 2, 3, 4$ are almost indistinguishable.

Recall that the periodicity measures are defined in terms of ε -simplifications of the expression profiles. The concept of an ε -simplification was introduced in [61], where the main result is a construction for functions on 2-manifolds. As described in this section, existence is obvious for functions on a 1-manifold. The situation is much less understood for functions on a 3-manifold. The question of the stability of the total persistence for Lipschitz functions was considered in [36]. Similar to the degree- q Wasserstein distance between diagrams studied in Section VIII.2, the difference between degree- q total persistences goes to zero as the functions approach each other for some values of q and not for others. For both concepts, the qualitative change happens at a value of q that depends on the dimension of the manifold on which the Lipschitz functions are defined.

IX.2 Elevation for Protein Docking

In this section, we express the protrusions and cavities of a surface using a real-valued function whose design is motivated by the 3-dimensional shape matching problem central to the molecular basis of life.

Background. According to the central dogma of biology, strands of DNA are transcribed to pieces of RNA, which are then translated into proteins. Transcription works by complementarity, while translation is more involved, going from an alphabet of four nucleotides to one of twenty amino acids. Proteins are made of strings of amino acids. These strings are highly variable in order and length. In principle, this suggests an astronomical number of different possible proteins, but nature apparently uses only a tiny fraction of perhaps a few hundred thousand types. Once a protein has been formed, it folds into a characteristic shape. This shape determines its function, that is, how the protein acts within its environment and, in particular, how it interacts with and binds to other proteins.

The interaction between proteins is one of the most fundamental processes in biology and holds the key to how biological systems work. Cells send signals to one another and build machines that perform the many tasks that make life possible. To understand these and other processes, it would be wonderful if we could predict which proteins interact with which other proteins simply by knowing their shapes and the forces exerted by their atoms. This is the *protein docking problem*, the computational prediction of protein interaction. However, this prediction has proven notoriously difficult. There is significant debate in the biochemistry community about the relative importance of the geometry (shape) and the physics (forces), but

it is clear that both are involved. It stands to reason that the relative importance of the geometry increases with the size of the involved molecules. But proteins flex, so geometry alone cannot predict the matching of undocked proteins. Nevertheless, geometric analysis is the first step.

Technology. The starting point for most docking efforts is geometric structures of proteins and other molecules collected by the biochemical community. The Protein Data Bank is an archive that contains information about experimentally determined geometric structures of proteins and nucleic acids. Data comes in the form of 3-dimensional atomic coordinates labeled by atom type and other descriptors. Data is determined primarily using two technologies, x-ray crystallography and nuclear magnetic resonance. For the former, biochemists crystallize the molecule and image the crystallized arrangement with a beam of x-rays that scatter in a variety of directions. From the angles and intensities of the scattered rays, a 3-dimensional picture of the density of electrons is produced. This density then allows the estimation of the positions of atoms in the crystal, as well as their chemical bonds. In contrast, nuclear magnetic resonance aligns nuclei with a magnetic field and perturbs this alignment with an orthogonal field. The response to the perturbation is then used to estimate the location of individual atoms.

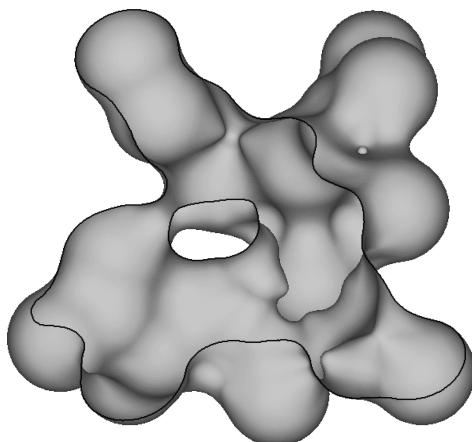


Figure IX.4: Cut-away view of a molecular skin surface.

Protein surfaces. Given atomic coordinates, we are interested in features on the surface of a protein that suggest binding configurations. The first problem is to define what we mean by the protein surface. Individual atoms attract and repel one another in several ways. The strong forces that hold the molecule together are chemical bonds and electrostatic interaction of ions. A weaker set of forces, known collectively as the *van der Waals force*, is strongly repulsive at short distance,

attractive at medium distance, and negligible at large distance. It favors a fixed distance along a large patch of contact, which is the reason why geometry plays a role in the interaction. To model this contact, we place a small sphere, called the *van der Waals sphere*, around the center of each atom. To define a surface, we keep the spheres fixed while rolling a ball about the configuration, always touching but never overlapping any of the spheres. The radius of the ball is chosen to approximate that of a water molecule. As the ball rolls around, it traces out the *molecular surface*, which is made up of sphere and torus patches. This is roughly how the surrounding water experiences the protein. Except for the occasional sharp edge formed by intersecting blending surfaces, the patches meet to form a continuous bundle of normal vectors. If continuity is important everywhere, we may alternatively use the *molecular skin*, which consists of sphere and hyperboloid patches; see Figure IX.4. However a protein surface is defined, we look for protrusions and cavities that might line up when two proteins interact. The mathematical tool we use to do this is called elevation, and it can be defined for curves in the plane or surfaces in 3-dimensional space. Although our primary application is to surfaces, we simplify the discussion by restricting ourselves to curves.

Curves in the plane. Suppose $\mathbb{M} \subseteq \mathbb{R}^2$ is the image of a smooth embedding of the circle. Define $F : \mathbb{M} \times \mathbb{S}^1 \rightarrow \mathbb{R}$ by mapping each point $x \in \mathbb{M}$ and each $u \in \mathbb{S}^1$ to the height of x in the direction u , that is, $F(x, u) = \langle x, u \rangle$. Fixing a direction, we get $f_u : \mathbb{M} \rightarrow \mathbb{R}$ defined by $f_u(x) = F(x, u)$, the height function in the direction u . We are interested in conditions for which this height function is Morse. Recalling the definition in Section VI.1, we note that f_u may fail to be Morse for two reasons, namely because it contains a degenerate critical point or it has two critical points sharing the same height value.

A simple degenerate critical point is modeled by the family of functions $g_s(t) = t^3 + st$. For $s < 0$, we have two critical points, one a local maximum and the other a local minimum. For $s = 0$, we have a degenerate critical point at $t = 0$, and for $s > 0$, we have no critical points. As s goes from negative to positive, the pair of critical points cancel each other. We call this event a *cancellation*, or an *anti-cancellation* if we go in the other direction. Similarly, we can use a parametrized fourth degree polynomial to model an *interchange* at which two critical points momentarily share the function value. In our case, varying the parameter, s , corresponds to moving the direction such that the critical points slide on the curve. A cancellation occurs when two critical points collide, which happens at an inflection point. This motivates us to assume that the curve has only a finite number of inflection points and only a finite number of lines that are tangent at two or more points. It follows that there are only a finite number of directions for which f_u is not Morse. Equivalently, the 1-parameter family of height functions passes through only a finite number of cancellations, anti-cancellations, and interchanges.

Elevation function. When f_u is a Morse function, we can use extended persistence to pair up its critical points. These are the points for which u is normal to the curve, and we associate to each the persistence of the pair to which it belongs,

calling this real number the *elevation* of the point. By the Persistence Symmetry Theorem of Section VII.3, the pairing is the same if we substitute $-u$ for u . It follows that f_u and $f_{-u} = -f_u$ define the same elevation values for the same points. In other words, elevation depends only on the normal line to the curve. Since every point of \mathbb{M} has a unique normal line, this defines the *elevation function*, $E : \mathbb{M} \rightarrow \mathbb{R}$, except at points at which the height functions are not Morse. We take limits to define E also at these exceptional points.

Recall that a cancellation happens at an inflection point, x , for which we set $E(x) = 0$. Indeed, it is easy to see that the two critical points are paired right before they meet and cancel each other at x . The limit is the same from both sides, namely zero, which implies that E is continuous at x . At the points involved in an interchange, however, we may have different left and right limits and thus an ambiguity of how to define E . This is illustrated in Figure IX.5. Here, the point x

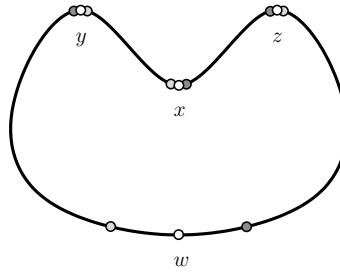


Figure IX.5: The four white points share the same normal direction, as do the four light shaded and the four dark shaded points.

would be paired with either y or z . In fact, if we rotate the vertical, upward directed vector u slightly to the right, the critical point near x is paired with the critical point near y . In contrast, if we rotate u slightly to the left, the critical point near x is paired with the critical point near z . Thus, there is a jump from y to z in moving the normal from right to left. The elevation near x varies continuously, so we can simply set $E(x)$ equal to the absolute height difference between x and y , which is the same as the absolute height difference between x and z . But the left and right limits at y are different, jumping from the absolute height difference between y and x to that between y and w . We get the same two different limits and the opposite jump at z . Continuity can therefore be obtained through surgery. Specifically, we cut \mathbb{M} at y and at z and we glue the four ends in pairs to get a new curve on which E is continuous at the cut points. Of course, the new curve is no longer embedded in the plane. In this particular case, the new curve consists of two components, a long loop that contains w and a short loop that contains x . If we perform surgery at all such discontinuities, we get a curve, \mathbb{N} , on which E is everywhere continuous.

Elevation maxima. Our interest in surgery is merely a means to an end, namely the determination of the interesting features of the curve. Call a point $x \in \mathbb{N}$ a *local maximum* of E if it has an open neighborhood such that $E(y) \leq E(x)$ for all y

in this neighborhood. For convenience, we assume the smooth curve is *generic* by which we mean

- (i) it has only finitely many height functions that are not Morse;
- (ii) its elevation function has only a finite number of local maxima.

Condition (i) has been discussed earlier, where it was used to define the elevation function. Condition (ii) prohibits curves of (locally) constant width, such as for example the circle. Consider the curve in Figure IX.5 as an example. Its elevation function has six local maxima, x and w , the two cut points formed by gluing the four ends obtained by cutting the curve at y and at z , as well as the leftmost point, p , and the rightmost point, q , of the curve (both not shown).

The local maxima come in pairs, by construction. For example, p and q form a pair, and $E(p) = E(q)$ is the Euclidean distance between p and q . Since neither point is a cut point, we call this pair a *one-legged elevation maximum*. We note that having p and q paired by extended persistence is necessary to form an elevation maximum but it is not sufficient. In the one-legged case, the line connecting p and q must be in the direction of the normal vector, and the curvature at p and at q must be such that a small rotation does not increase the local width. A second pair is formed by x and the cut point produced by gluing the right end at y to the left end at z . We call this a *two-legged elevation maximum* because we have two legs connecting x to y and to z on the original curve. Again, having x paired to y and to z by extended persistence is necessary but not sufficient to form a two-legged elevation maximum. We also need the property that the orthogonal projection of x onto the line of y and z falls between the points y and z . Finally, we have a third pair formed by w and the cut point produced by gluing the left end at y to the right end at z . This is another two-legged elevation maximum.

Piecewise linear curves. To design an algorithm for computing the elevation maxima of a curve, we face the usual dilemma that input is never smooth. Instead, we assume a simple, closed polygon with vertices x_0, x_1, \dots, x_{n-1} and edges e_i connecting x_i to x_{i+1} , for $0 \leq i < n$ where we take indices modulo n . We assume the polygon is *generic*, by which we mean that no two of the $\binom{n}{2}$ lines connecting the n vertices are parallel or orthogonal to each other. We may think of this polygon as approximating a smoothly embedded curve and this way get an idea of what the elevation maxima ought to be. Alternatively, we may approximate the polygon by a generic, smooth curve and obtain the definitions by limit considerations. This is what we do next.

Let P be the subset of \mathbb{R}^2 whose boundary is the polygon. We write P^ε for the set of points at distance at most ε from P and $(P^\varepsilon)^{-\delta}$ for the set of points of P^ε at distance more than δ from the complement. For $\delta = \frac{\varepsilon}{2}$ and sufficiently small $\varepsilon > 0$, the boundary of $(P^\varepsilon)^{-\delta}$ alternates between an arc on a circle with center x_i and radius δ and a straight line segment parallel to e_i . Finally, we replace the straight line segment by circular arcs of ever so small, positive curvature, κ , calling them the *chords* connecting the circular arcs around the vertices. While the resulting curve

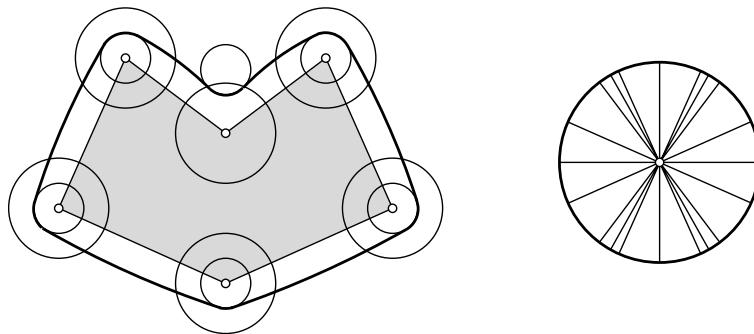


Figure IX.6: Turning a polygon into a generic curve. To simplify the drawing, we have chosen to ignore the requirements for a generic polygon. The circles at the vertices have radius ε and $\delta = \frac{\varepsilon}{2}$; they illustrate the construction of $(P^\varepsilon)^{-\delta}$. On the right, we see the circle of directions decomposed into arcs of constant height ordering.

is not smooth, its normal bundle is continuous and it satisfies the requirements of a generic curve, which suffices to define the elevation function; see Figure IX.6. For sufficiently small $\varepsilon > \kappa > 0$, the points on a chord are all paired with points on the same circular arc. In contrast, points on a circular arc may be paired with points of more than one chord or arc.

Algorithm. We begin by recalling the Extended Persistence Algorithm applied to a height function, f_u , defined on the polygon. Assume the vertices have distinct height values and they are relabeled such that $f_u(y_0) < f_u(y_1) < \dots < f_u(y_{n-1})$. Then each vertex can be unambiguously classified as a minimum, a regular vertex, or a maximum. Since minima and maxima alternate along the polygon, we have the same number of each, and the algorithm outputs a pairing (a perfect matching) between the two collections. The global minimum, y_0 , is necessarily paired with the global maximum, y_{n-1} . The other pairs depend on the sequence of sublevel sets or, equivalently, the lower star filtration of f_u , as explained in Chapter VII. We note that the pairing is the same for other piecewise linear functions for which the ordering of the vertices is the same. It thus suffices to run the Extended Persistence Algorithm for $\binom{n}{2}$ height functions, one per antipodal pair of arcs defined by the vertex pairs; see Figure IX.6. Skipping a few details, we note that this can be done in time at most some constant times n^3 .

We now discuss how the $\binom{n}{2}$ pairings are used to extract all elevation maxima. As mentioned earlier, there are two types. We first discuss the one-legged elevation maxima. Let x_i, x_j be two vertices and $u = (x_j - x_i)/\|x_j - x_i\|$ the direction they define. Then x_i and x_j form a one-legged elevation maximum iff x_i and x_j are paired by the algorithm applied to f_u . In the piecewise linear case, the condition on the curvature at the two points is void. We second discuss two-legged elevation maxima. Let x_i, x_j, x_k be three vertices and let u be normal to $x_k - x_j$. Suppose

furthermore that x_j and x_k lie on opposite sides of the line with direction u that passes through x_i . Let u_- and u_+ be directions sufficiently close to and on opposite sides of u . Then x_i and x_j, x_k form a two-legged elevation maximum iff x_i and x_j are paired for f_{u_-} and x_i and x_k are paired for f_{u_+} . In summary, the $\binom{n}{2}$ runs of the Extended Persistence Algorithm provide all the information we need to identify the elevation maxima of the polygon.

Notes. The difficulty of predicting the binding between proteins of known geometric structure combined with the importance of this question has lead to a community organized competition [87]. Using yet unpublished geometric structures determined by x-ray crystallography [52] or by nuclear magnetic resonance [159], the participants are asked to submit their best predictions, which are then compared to the observed configuration. The idea of using protrusions and cavities of protein surfaces to predict binding configurations goes back to Connolly [40]. He represents the protein by its molecular surface, which decomposes \mathbb{R}^3 into the *inside* and the *outside*, two 3-manifolds with disjoint interiors and common boundary, namely the molecular surface. Fixing a radius, $r > 0$, he places the center of a sphere with radius r at every point x of the surface and assigns to x the fraction of the sphere contained in the inside. As shown in [29], the limit of this function, as r approaches zero, is related to the mean curvature function of the surface. This should be contrasted to the relationship between the total mean curvature and the elevation that follows from integral geometric considerations described in [33]; see also Section VIII.3. As demonstrated in [152], the elevation maxima are useful in the coarse alignment of protein structures. This suggests we use elevation as a first pass toward predicting a binding configuration and refine the resulting alignments with methods that incorporate detailed knowledge of the physical behavior of molecular systems [130].

While this section focuses on the simpler setting of a curve embedded in \mathbb{R}^2 , the important setting is of course that of a surface embedded in \mathbb{R}^3 . This is described in the original paper on the subject by Agarwal, Edelsbrunner, Harer, and Wang [3]. In going from curves to surfaces, the ideas remain the same but get technically more complicated. For smoothly embedded surfaces, the construction is based on Cerf theory [30], which is part of differential topology. Instead of two, we get four types of elevation maxima; see Figure IX.7. Except for the one-legged case, each

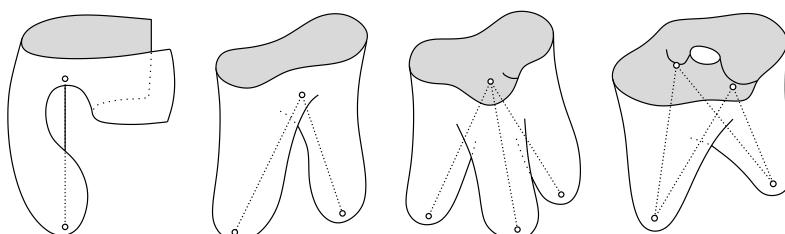


Figure IX.7: From left to right: a one-, two-, three-, and four-legged elevation maximum of a surface embedded in 3-dimensional space.

type arises at an ambiguity of the pairing of critical points produced by extended persistence. The algorithms are again for piecewise linear functions. Similar to the polygonal case, we construct all elevation maxima by running the Extended Persistence Algorithm for a finite collection of height functions. We refer to Section VII.2 for a fast implementation of the Persistence Algorithm for piecewise linear functions on a 2-manifold. Computing the extended persistence pairs is, however, more difficult and a similarly efficient algorithm requires sophisticated data structures [77].

IX.3 Persistence for Image Segmentation

In this section, we discuss image data and, in particular, the problem of segmenting the data into meaningful pieces. A popular approach is the watershed method, but it is sensitive to noise in the data, tending to overdo the segmentation. We show how to use persistent homology to cope with this difficulty.

Background. When we collect data about a physical phenomenon, we do so to varying degrees of resolution. Images are high-resolution data sets, representing shapes and scenes in great detail. A large part of biological and medical research, as well as medical practice, depends on technology that produces 2- and 3-dimensional images. But we can go beyond three dimensions, eg. with video sequences that unwind in time. There are also reasons for generating images synthetically, using methods such as Fourier transforms, for finding symmetries and for other purposes. The high resolution of image data suggests we think of it as a continuous object and apply methods from continuous rather than discrete mathematics for its analysis. By its nature, an image contains more than the desired information. Therefore the first task is often the extraction of interesting features. Capturing and describing these features is the province of image analysis. It includes tasks such as denoising, segmentation, registration, comparison, and more.

Technology. The last decades have witnessed a revolutionary change in how science is practiced, and this change is fueled by ever improving ways of acquiring data. Using new technology, we are able to collect high-resolution data on physical events that have traditionally been beyond our reach. Examples are sensor networks monitoring environments and microarrays measuring the expression of the entire genomes. These are relatively recent technologies for which we can expect rapid improvements in the accuracy and volume of collected data. More traditional imaging technologies generate 2- and 3-dimensional arrays of measurements.

Microscopy. This is an umbrella under which we distinguish different technologies depending on the medium used to generate the image. Optical microscopy involves light diffracted from an object passing through a lens to allow for a magnified view. Similarly, electron microscopy measures the diffraction of electromagnetic radiation by an object. In contrast, scanning probe microscopy, as the name suggests, measures the interaction of a probe with an object. In each case, the output is a

2-dimensional array of tiny squares, referred to as *pixels*, short for picture elements. Being 2-dimensional, the number of pixels is usually not more than perhaps a few million, which is easy to manage with current computer storage technology.

Magnetic resonance imaging. The principles are the same as for nuclear magnetic resonance, but the N-word has been dropped for medical applications. Here we use a magnetic field to align the nuclei of hydrogen atoms in water. Radio frequency fields are then used to systematically alter the alignment, which creates the signal detected by the scanner. This technology is widely employed in radiology to study the internal structure of the human body. The output is a 3-dimensional array of tiny cubes, referred to as *voxels*, short for volume elements. Being 3-dimensional, the voxel array provides a lot more information than a pixel array generated with microscopy. This wealth comes with a cost, namely the added difficulty of managing and analyzing such a large amount of data.

Segmentation. Once an image is acquired, it becomes a mathematical object that we can study. In particular, a 2-dimensional image is a function that is piecewise constant on a rectangular configuration of pixels; see Figure IX.8. Inside a pixel, the function is constant and measures intensity or gray value, and similarly for color pictures, except that we get three separate images, one for red, one for green, and one for blue. Given an image, the *segmentation problem* looks to identify

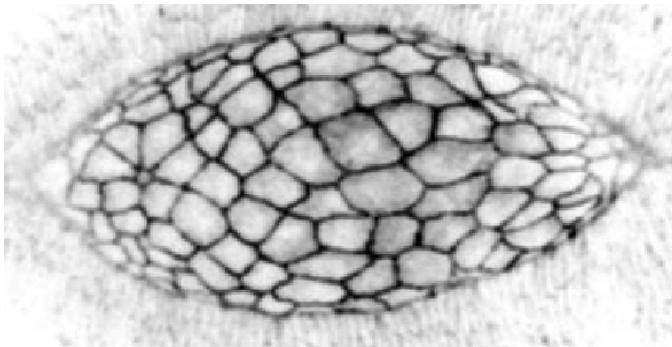


Figure IX.8: Confocal microscopy image of a cross-section of the cells in a Drosophila embryo during the developmental process known as dorsal closure [image courtesy of Daniel P. Kiehart and Adrienne Wells, Biology, Duke University].

regions of interest from these values. In Figure IX.8, these regions would be the cells imaged by microscopy. The general problem is hopelessly ill-defined but of major importance. As a result, the state of the art in the field is at best imperfect. Every type of image provides a different set of challenges, motivating a variety of approaches to the problem. Specifying global or adaptive thresholds to define the regions is a good first step. Methods of mathematical morphology can then be used to refine the result. Deformable models or level set methods solve differential equations to shrink-wrap a region with a curve or a surface. Region growing

and region splitting methods seek to improve segmentations by locally improving a quality assessment. In this section, we consider yet another approach, one that fits nicely into the framework of computational topology.

Watershed intuition. We begin with an intuitive description of the method. Let us treat an image as a continuous function defined on a region of the plane, usually a rectangle, although for this description we assume it is defined on all of \mathbb{R}^2 . Think of the graph of this function, a surface in \mathbb{R}^3 , which we imagine permeable, with soil below and air above. Now suppose it rains and the water level rises everywhere on the plane. As is common on planet Earth, we call *land* only the part of the surface above the water level. As the level rises, we see the land shrinking and its topology changing. When the water reaches a local minimum, a lake forms and grows as the level rises. When the water reaches a saddle point, two lakes merge into a single lake or an island separates from the mainland. When the water reaches a local maximum, the corresponding island has completely submerged under water.

We can keep the water from overflowing by building watershed lines as we pass saddles. These are the curves that separate lakes where the water meets as it rises. Mathematically, they form the unstable or ascending 1-manifolds corresponding to the saddles. They prevent the lakes from merging and form roads (anti-canals) between the islands and the mainland, thus maintaining the lakes as topological disks throughout the process.

The Watershed Algorithm. To formalize this process, we construct a piecewise linear function that represents a given image. For each pixel, we have a vertex at its center and we connect the vertices to form a triangulation. It is convenient to compactify by adding a dummy vertex to get a triangulation of S^2 . Specifying a value at each vertex, we get a function by piecewise linear interpolation; see Section III.1. More generally, we can begin with a triangulated 2-manifold and a piecewise linear function, $f : M \rightarrow \mathbb{R}$. Recall that in Section VI.2, we constructed a complex whose vertices were the maxima, edges were the unstable 1-manifolds, and regions were the unstable 2-manifolds of the function. We now give an algorithm that constructs an approximation of the unstable manifolds. It is convenient to assume that the vertices have distinct function values and they are indexed such that $f(x_1) < f(x_2) < \dots < f(x_n)$. We recall from Section VI.3 that each vertex is classified as regular or critical by looking at the values of f in its link. Call the part of the link spanned by vertices whose function value exceeds $f(x_i)$ the *upper link* of x_i . Then x_i is a minimum if the upper link is the entire link, a maximum if the upper link is empty, and a k -fold saddle if the upper link consists of $k + 1$ components, each a path or an isolated vertex. The vertex is regular if it is a 0-fold saddle, that is, if it has a non-empty, connected, upper link that does not exhaust the entire link.

We process the vertices from lowest to highest. Specifically, we run a loop from $i = 1$ to $i = n$ and distinguish between the different types of vertices. Initially, all simplices in the triangulation are unmarked.

CASE 1: x_i is a saddle. We mark x_i together with the edges that connect the saddle to the highest vertex in each component of the upper link.

CASE 2: x_i is regular. If x_i has an incoming marked edge, then we mark x_i together with the edge to the highest vertex in the upper link.

CASE 3: x_i is a maximum. We mark x_i .

In the end, we have $k + 1$ paths running upward from a k -fold saddle. Sometimes these paths merge, but then they continue together until they reach a maximum. The number of paths ending at a maximum varies depending on the surrounding configuration of minima and saddles. It is even possible that a maximum is isolated, without a path ending at the vertex. But this can only happen if the manifold is a sphere and f has one minimum, one maximum, and no saddles. It is not difficult to see that the paths consisting of the marked edges and vertices cut the 2-manifold into open disks, one for each minimum. This is also true if we have no saddles and therefore no marked edges. The open disk is then the sphere minus the maximum, which is marked by the algorithm.

Cleaning up. The Watershed Algorithm is widely used, but it tends to overdo the segmentation, creating too many regions and identifying small noise in the image rather than just the desired features. For this reason, there is always a clean-up step, sometimes done systematically and sometimes in an ad hoc way or even manually. This is illustrated in Figure IX.9, where the goal of the segmentation is to identify the location and the shape of cells of a fly embryo. To appreciate the importance of a reliable and consistent segmentation, we note that the image shows a cross-section of the embryo and similar images are taken at other cross-sections. Furthermore, the images are taken in a time series. After segmenting each cross-section, the task is to connect the results to reconstruct the 3-dimensional cells and then the cells to reconstruct the motion. Finally, the details of the motion are used as cues to hypothesize the forces that drive the motion. We can clearly see that the segmentation in Figure IX.9 at the top is too fine to capture individual cells. We need to simplify the segmentation by distinguishing more from less important separations. Persistence gives us just the tool we need for this.

Simplification. We begin by computing the persistence of the minima, saddles, and maxima, which can be done during the same bottom-up sweep that constructs the segmentation. We get infinite persistence for the critical vertices, giving birth to essential homology classes, and finite, positive persistence for all other critical vertices. For simplicity, assume that all saddles are 1-fold and thus get assigned a unique persistence value, the absolute height difference to the paired minimum or maximum.

We simplify the segmentation in the order of increasing persistence. Assuming the absolute height differences of the vertex pairs computed by the Persistence Algorithm are distinct, the pair with minimum persistence is unique. Consider first the case in which this pair consists of a minimum, x , and a saddle, y . Passing y during the sweep merges the component started at x with another component

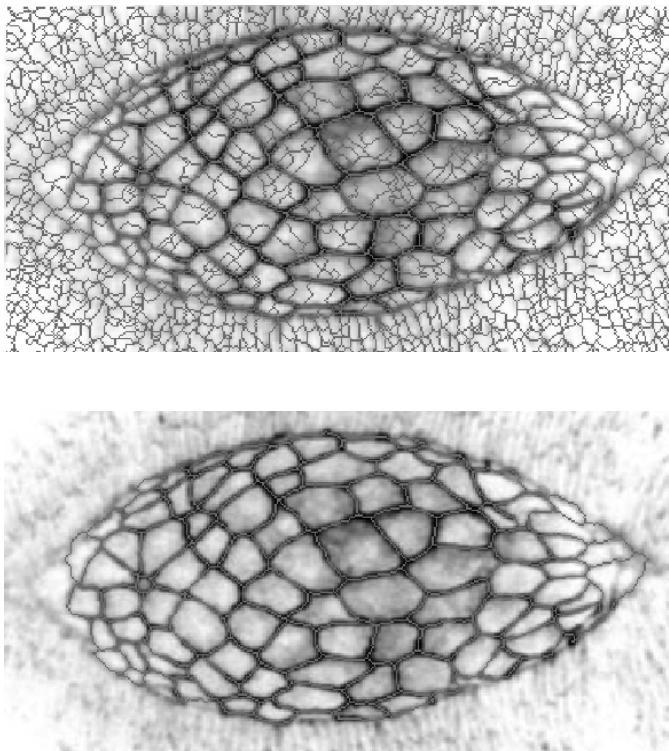


Figure IX.9: Top: the initial watershed segmentation before clean-up. Bottom: the result of simplifying the segmentation using persistence.

started earlier. By minimality of persistence, the watershed line started at y is part of the boundary of the region of x . Therefore, we can simplify by removing this line, which we do in two passes, both beginning at y . Unmarking edges and vertices in sequence, we stop each pass when we reach a maximum or the line merges with another. The result is that the two regions separated by y are now joined, and since x was the minimum of one of these, the other minimum represents the merged region. Consider second the case in which the minimum persistence pair consists of a saddle, y , and a maximum, z . By minimality of persistence, the watershed line started at y ends at z on one side and at a higher maximum on the other side. If this is the only watershed line ending at z , we remove it by unmarking its vertices and edges in sequence, beginning at z . Otherwise, we let the watershed line be, except that we think of it as an extension of the other lines ending at z .

We note that the change effectively treats the critical vertices of the minimum persistence pair as if they were regular. The rest of the persistence pairing remains unchanged. We can therefore proceed to the next lowest persistence pair and continue until we exceed a pre-chosen threshold. Applying this strategy to the segmentation in Figure IX.9, top, we get the segmentation shown at the bottom.

Notes. Images are generated by a plethora of technology, including microscopy [42] and magnetic resonance [81]. The 2-dimensional images in this section are from work on the dorsal closure in fly embryos [91]. Algorithms for processing images are described in the image analysis literature [135]. Many problems in this area are of a topological nature [93], including the segmentation of images into regions of interest. The Watershed Algorithm for segmentation goes back to the early eighties of last century [18]; see [128] for a survey of the general literature on the topic. Because of the importance and the large amount of medical data, the 3-dimensional version of the algorithm is of particular interest. We refer to [131] for the description of the method for magnetic resonance images using a diffusion filter to cope with the endemic over-segmentation. The algorithm for 3-dimensional images is similar to but more complicated than for 2-dimensional images. From Morse theory, we know that we have four types of simple critical vertices: minima, index-1 saddles, index-2 saddles, and maxima. We get a 3-dimensional cell for each minimum, a surface for each index-1 saddle, a curve for each index-2 saddle, and a vertex for each maximum. Together, they form a complex akin to the unstable manifolds discussed in Chapter VI. Persistence pairs minima with index-1 saddles, index-1 with index-2 saddles, and index-2 saddles with maxima. The simplification can again be done in the order of increasing persistence, but this is now more complicated than for 2-manifolds.

IX.4 Homology for Root Architectures

In this section, we look at the problem of recovering the structure of a plant root from photographic images. We combine standard image processing techniques with homology computations to capture interesting traits, such as the branching pattern and the distribution in space.

Background. Plant biologists understand much more about how plants grow and develop above the ground than underground. Yet, the root is every bit as important in how a plant responds to environmental variation and how it adapts to soil and moisture conditions. Learning about root architecture beyond what we know today is necessary before we can begin to understand the connection between phenotype and genotype in root development. The genotype is studied in a variety of biological experiments, many involving microarrays used to measure the expression of an entire collection of genes. To characterize the phenotype, we need an accurate set of measurements, preferably obtained without moving or damaging the plant. This way we can repeat the measurements during development, while the root makes decisions about where and when to grow.

We focus here on topological features of a plant root, in particular on a decomposition into tips, forks, and branches and on a characterization of space utilization. At a fork, a growing root either divides into two or a lateral root emerges. If we remove the fork, the rest of the root is branches, and we call the end of a branch that is not a fork a tip. Plant roots grow from the tip, so the number and location

of these is of importance to biologists. Note, however, that the location of the forks and tips along the root says little about the way the root distributes itself in the soil. To study this distribution, we consider the complement of the root embedded in space and measure its connectivity using persistent homology.

Technology. For simple and rather obvious reasons, studying the architecture of a plant root is a difficult undertaking. We can dig up the plant to measure traits like length, branching, and more, but there are limitations to this approach. The first is that removing a plant from the ground seriously disrupts its growth pattern and may even kill it. A second is a lack of information about space utilization. We would like to find out how the root distributes itself in the soil and how its growth varies in response to a variety of stimuli, such as soil nutrient abundance and distribution, other forms of environmental stress, and the availability of water.

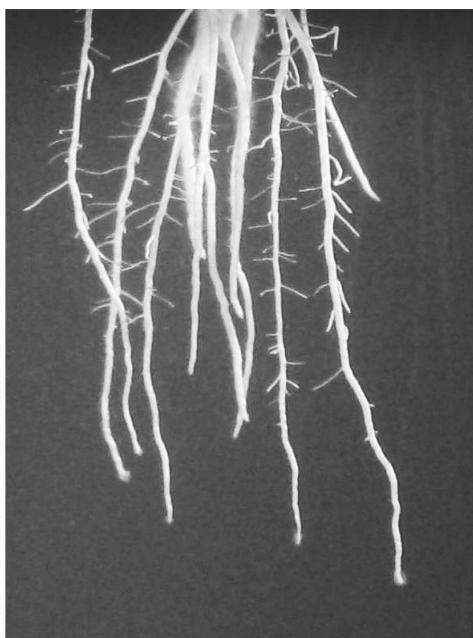


Figure IX.10: Rice root system growing in gel [image courtesy of Philip N. Benfey and Anjali Iyer-Pascuzzi, Biology, Duke University].

To cope with these difficulties, we need a new medium to grow the plant and a way to image the root. We know of two solutions: growing the plants in styrofoam containers and taking x-ray images, and growing them in transparent gel and taking photographs. The styrofoam provides a fairly realistic medium for the plant and both nutrition and water conditions are easy to vary without disturbing the plant. The disadvantage is the need to vacuum the container to remove as much water as possible before taking images. Water is an issue because x-rays refract when they pass through water, so any moisture left in the container shows up as noise in the

image. The problem is severe since vacuuming disturbs the plant while moisture compromises the images. The transparent gel provides a nutrient mix that is less realistic than that of the styrofoam but much more than that of a hydroponic system, for example. The main advantage is the ease with which we can take photographic images; see Figure IX.10. Placing the gel together with the root inside a glass container, we can take photographs 360 degrees around. The task thus reduces to extracting the desired information from a 2-parameter sequence of images, going around the root and taking the photographs over a period of a few weeks. Each image is a 2-dimensional array of pixels. We discuss the extraction of features directly from these images as well as attempts to reconstruct a 3-dimensional image from the sequence of 2-dimensional images.

Tips, forks, and branches. Suppose first that we are working with a single photographic image, that is, a projection of our root to the plane represented by an array of pixels, p , with intensities, $f(p)$. Specifying a threshold, θ , we decompose the image into *foreground*, the union \mathbb{Y} of all pixels with intensity $f(p) \geq \theta$, and *background*, the union \mathbb{X} of pixels with intensity $f(p) < \theta$. We recall that each pixel is a closed square so that foreground and background are both closed and intersect in their common boundary, which is 1-dimensional.

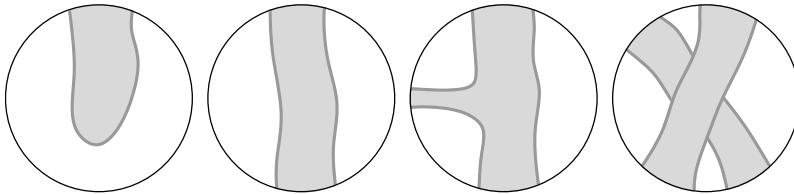


Figure IX.11: From left to right: schematic local pictures of a tip, a branch, a fork, and a crossing.

Assuming the foreground represents the root, it consists of streams of pixels forming roads that fork and cross and eventually end. While the roads vary in their thickness, we think of them as forming a 1-dimensional graph, with nodes connected by arcs. We call a degree-1 node a *tip*, a degree-3 node a *fork*, and an arc between two such nodes a *branch*. We illustrate the definitions by showing typical local pictures in Figure IX.11. There are also degree-4 intersections, but since roots rarely sprout two lateral branches at the same location, we assume these are artifacts of occlusion and represent them as crossings between branches rather than nodes in the graph. Of course, there can be more complicated situations caused by accumulated occlusion. These are better dealt with in three dimensions, and we ignore them for the time being.

Persistent local homology. As suggested by Figure IX.11, we look at the homology of the foreground and the background within small circular windows to classify a pixel to belong to a tip, a fork, or a branch. Fixing a point $x \in \mathbb{R}^2$ and

a real number $r \geq 0$, we write $B(r) = B_x(r)$ for the closed disk with center x and radius r . The foreground and background within this window are $\mathbb{Y}(r) = \mathbb{Y} \cap B(r)$ and $\mathbb{X}(r) = \mathbb{X} \cap B(r)$. We are interested in the first homology of the foreground within the window relative to its boundary on the circle, $H_1(\mathbb{Y}(r), \text{bd } B(r))$. Assuming the generic case in which the circle intersects the boundary of \mathbb{Y} transversally, we replace the boundary on the circle by the boundary $\mathbb{Y}(r)$ shares with $\mathbb{X}(r)$. Since $\mathbb{Y}(r)$ is in the plane, it is easy to see that this gives an isomorphic first homology group; compare with Exercise 7 at the end of this chapter. We thus get

$$\begin{aligned} H_1(\mathbb{Y}(r), \text{bd } B(r)) &\simeq H_1(\mathbb{Y}(r), \mathbb{Y}(r) \cap \mathbb{X}(r)) \\ &\simeq H_1(B(r), \mathbb{X}(r)) \\ &\simeq \tilde{H}_0(\mathbb{X}(r)), \end{aligned}$$

where we get the second line by excision and the third line using the exact sequence of the pair $(B(r), \mathbb{X}(r))$. Instead of looking at the first relative homology group, we can therefore use the zeroth absolute homology group of the local complement, $\mathbb{X}(r)$, to distinguish between the different types of neighborhoods; see Table IX.1. Using persistence, we eliminate the dependence on the choice of r . Specifically, we

	tip	branch	fork	crossing
rank $H_0(\mathbb{X}(r))$	1	2	3	4

Table IX.1: The rank of the zeroth homology groups of the neighborhoods of a pixel inside a tip, a branch, a fork, and a crossing.

increase r from zero to infinity and consider the zeroth persistence diagram of the local background pictures, $\mathbb{X}(r) \subseteq \mathbb{X}(s)$ for $0 \leq r \leq s < \infty$. If x is part of a tip, we see the following typical behavior as we grow the window.

Tip. For very small r , $\mathbb{X}(r)$ will be empty. Its first component will be born when r reaches the distance from x to \mathbb{X} . This might be the only event for a while, but more likely we will see births and deaths of components in quick succession. However, these extra components correspond to points in the diagram whose persistence is negligible. Of course, once r gets large, all kinds of things may happen. In summary, we see only one birth that happens for small r and whose persistence is not negligible.

Similarly, for a branch we see two births for small r with larger than negligible persistence, for a fork we see three such births, and for a crossing we see four. As one can imagine, using persistence instead of a fixed radius greatly increases the number of pixels that can be correctly classified, but it is still a long shot from classifying all pixels. Ambiguities arise for a variety of reasons, including spurious foreground and background components, thicker than expected branches, and other root portions reaching into the local window. We can conceive of heuristics coping with these difficulties, but ultimately we need to face the fact that the problem as described does not admit a perfect solution.

3-dimensional reconstruction. Important information about the root, including estimates for the number of tips and forks, can be computed directly from the 2-dimensional images. However, to learn how the root distributes itself to explore the soil requires a reconstruction as a subset of 3-dimensional space. We describe

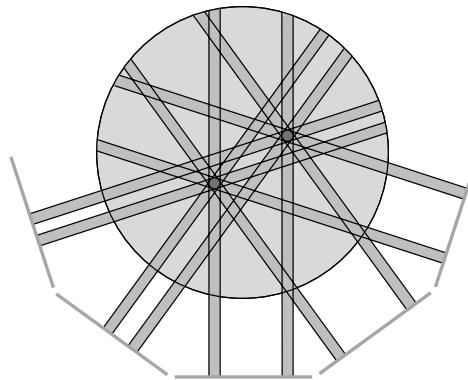


Figure IX.12: Schematic cross-section of a root growing inside a cylindrical container. We reconstruct the two shaded spots from their images in the five projections.

here an image processing approach to this problem. It starts with a cubic block of voxels from which the algorithm sculpts the root. We assume a small, but not too small, number of photographic images taken of the same root from different directions and at about the same time. For each image, we know the position of the camera and the direction of the projection. As before, we use a threshold to decompose an image into foreground and background. If the projection of a voxel into the plane of the image lands inside the background, then the voxel cannot be part of the root. As illustrated in Figure IX.12, the Space Carving Algorithm combines the information gleaned from all 2-dimensional images and this way arrives at a first approximation of the 3-dimensional structure.

The quality of the reconstruction depends on the number and resolution of the images, the calibration of the cameras, and other factors. As before, we can make amendments to the algorithm to improve the quality, such as estimating probabilities for a voxel to belong to the root or using prior knowledge about the structure of the root. While we can perhaps reach acceptable results, we should keep in mind that perfection is at best reachable in the limit of our improvement efforts.

Utilization of space. Suppose now that we have reconstructed the 3-dimensional structure of the root. Reusing the 2-dimensional notation, we write $\mathbb{Y} \subseteq \mathbb{R}^3$ for the space we use to represent the root. We may revisit the decomposition of the root into tips, forks, and branches using local homology within spherical balls. Since crossings no longer confuse the picture, we can expect a higher success rate in the classification. We can also address the global question of how the root distributes itself in space. For this purpose, we introduce the Euclidean distance function,

$f : \mathbb{R}^3 \rightarrow \mathbb{R}$, defined by mapping every point $x \in \mathbb{R}^3$ to its distance from \mathbb{Y} , that is, $f(x) = \inf_{y \in \mathbb{Y}} \|x - y\|$. Note that $\mathbb{Y} = f^{-1}(0)$. The sublevel sets of f form a nested sequence of spaces, $\mathbb{Y}_r = f^{-1}[0, r]$. The corresponding sequence of reduced homology groups,

$$0 \rightarrow \tilde{\mathcal{H}}_p(\mathbb{Y}_0) \rightarrow \dots \rightarrow \tilde{\mathcal{H}}_p(\mathbb{Y}_r) \rightarrow \dots \rightarrow 0,$$

characterizes how thickening up the root fills space. As described in Chapter VII, we use the persistence diagram to characterize the main events in the filtration. The root is connected, so the zeroth diagram, $\text{Dgm}_0(f)$, should be empty. Any deviation from this ideal will have to be explained by failures to accurately reconstruct the root. There is more interesting information in the first and second diagrams. Consider for example a diffuse root system, that is, a root that distributes itself reasonably densely and more or less uniformly in the available space. Then \mathbb{Y}_r will have trivial first and second homology groups already for small values of r . Correspondingly, $\text{Dgm}_1(f)$ and $\text{Dgm}_2(f)$ will have no points of larger than negligible persistence. On the other hand, if the root has a tendency to grow around pieces of space, then this will express itself in voids and tunnels of \mathbb{Y}_r . Correspondingly, one of the two or both diagrams will have points with larger than negligible persistence. We note that this discussion neglects the possibility of a root that grows radially in a non-uniform manner but avoids the creation of tunnels and voids while exploring space. But we can detect such behavior by considering spherical cross-sections, for example.

There is more than one way we can compute the persistence diagrams of f . For example, we can grow \mathbb{Y}_r by successively adding voxels to the initial space, $\mathbb{Y} = \mathbb{Y}_0$. Alternatively, we may let $S \subseteq \mathbb{R}^3$ be the set of centers of the voxels constituting \mathbb{Y} . We then compute the Delaunay complex of S and the family of alpha complexes, as explained in Chapter III. The Stability Theorems of Chapter VIII imply that the diagrams we get from these and reasonable other methods are only a small bottleneck distance away from each other.

Notes. The study of plant roots has a long tradition in biology [26]. The project that provides the background for the discussions in this section targets agricultural plants, such as rice and maize. Local homology is a natural choice in the study of structural features such as forks and tips in roots. In mathematics, the concept makes its first appearance in Poincaré duality. Letting \mathbb{M} be a d -dimensional manifold without boundary and x a point of \mathbb{M} , we find that the relative homology group $\mathcal{H}_p(\mathbb{M}, \mathbb{M} - \{x\})$ is trivial for all dimensions $p \neq d$ and has rank one for $p = d$. Using integer coefficients, the latter group has two generators, and a choice of one is called an orientation of \mathbb{M} at x . Making consistent choices at all points, an element of $\mathcal{H}_d(\mathbb{M})$ is a fundamental class of \mathbb{M} if its image under the induced map to $\mathcal{H}_d(\mathbb{M}, \mathbb{M} - \{x\})$ is the chosen orientation. This class is used to define Poincaré duality; see Hatcher [82, Section 3.3]. The idea of using this construction in combination with persistent homology appears for the first time in [17]. Given a finite point set $S \subseteq \mathbb{R}^d$, the paper uses persistent versions of local homology towards reconstructing a stratified space that best approximates S . Basic tools in

this study are the persistent kernels and cokernels of maps from one filtration to another. The algorithms for these are similar to but more involved than those for ordinary persistence [37].

The problem of reconstructing shapes from sequences of images is studied in the general field of computer vision; see e.g. [72, 86]. The idea of carving out the shape from a block of voxels is due to Kutulakos and Seitz [98], but see also [100]. Given such a reconstruction, we can use standard algorithms for Delaunay complexes [54], alpha shapes [63], and persistent homology [60] to characterize the distribution of the root in space.

Exercises

The credit assignment reflects a subjective assessment of difficulty. A typical question can be answered using knowledge of the material combined with some thought and analysis.

1. **Antipodal functions** (two credits). A function $f : \mathbb{S}^1 \rightarrow \mathbb{R}$ is *antipodal* if $f(s) = f(-s)$ for all $s \in \mathbb{S}^1$. Equivalently, the function defined by $g(s) = f(s) - f(-s)$ is the zero function.
 - (i) Design a measure for quantifying the distance of a function from being antipodal.
 - (ii) Prove that your measure is stable or, alternatively, change your measure such that it is stable.
2. **Lipschitz function on the sphere** (three credits). Let d be a positive integer constant and $f : \mathbb{S}^d \rightarrow \mathbb{R}$ a Lipschitz function. Note that the d -dimensional volume of \mathbb{S}^d is bounded from above by a constant that depends on d .
 - (i) Prove that for $q > d$, the degree- q total persistence of f is bounded from above by a constant.
 - (ii) Use the result in (i) to show that for $q > d + 1$, the degree- q total persistence measuring Lipschitz functions on \mathbb{S}^d is stable.
3. **Fast pairing** (two credits). Let $f : \mathbb{S}^1 \rightarrow \mathbb{R}$ be a continuous, piecewise linear function specified by its values at the vertices of a triangulating n -gon.
 - (i) Assuming $f(x_i) \neq f(x_j)$ for all pairs of vertices $x_i \neq x_j$ of the n -gon, characterize the vertices that are paired by extended persistence.
 - (ii) Furthermore assuming the vertices are given in the order of increasing function value, show that the extended persistence pairing can be computed in time at most some constant times n .
4. **Inflection points and bitangent lines** (one credit). Let $\gamma : \mathbb{S}^1 \rightarrow \mathbb{R}$ be a smooth embedding of the circle in the plane. Suppose the curvature vanishes only at a finite number of points. Show that there are only a finite number of lines that are tangent to the curve at two or more points.

5. **Labeling regions** (two credits). Consider the Watershed Algorithm for segmenting a triangulated 2-manifold given in Section IX.3.
- Modify the algorithm so it labels the simplices in each region with the index of the generating maximum.
 - Define the *i-th region* as the union of the interiors of the simplices labeled *i* by the modified Watershed Algorithm. Prove that it is homeomorphic to an open disk.
6. **Ordering the pixels** (two credits). Let $n = 2^k$ and consider an n -by- n array of pixels p_i^j for $1 \leq i, j \leq n$. We define what it means to list the pixels in *Z-order*. For $k = 1$, we have four pixels which we arrange as $p_1^1, p_1^2, p_2^1, p_2^2$. For $k > 1$, we decompose the array into four equal blocks and list the upper left, the upper right, the lower left, the lower right blocks in this sequence and each in *Z-order*.
- Assume the pixels are listed in lexicographic ordering of their index pairs, (i, j) . Write an algorithm that rearranges the pixels in *Z-order*.
 - Write computer programs that translate back and forth between the row-column index pairs of a pixel and its index in *Z-order*.
7. **Isomorphic relative homology groups** (three credits). Let \mathbb{Y} be a d -manifold with boundary and let $\text{bd } \mathbb{Y} = \mathbb{A} \cup \mathbb{B}$ be a decomposition of the boundary into two $(d-1)$ -manifolds with common, $(d-2)$ -dimensional boundary and disjoint interiors. Prove that $H_{d-p}(\mathbb{Y}, \mathbb{A}) \simeq H_p(\mathbb{Y}, \mathbb{B})$ for all dimensions p .
8. **Distance function** (two credits). Let S be a finite set of points in \mathbb{R}^d and let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined by $f(x) = \min_{u \in S} \|x - u\|$. Recall from Chapter III that $\text{Alpha}(r)$ is the alpha complex defined by S and a radius $r \geq 0$.
- Let $r \leq s$ and consider the diagram of homology groups in which all four maps are induced by inclusion:
- $$\begin{array}{ccc} H_p(f^{-1}[0, r]) & \longrightarrow & H_p(f^{-1}[0, s]) \\ \uparrow & & \uparrow \\ H_p(\text{Alpha}(r)) & \longrightarrow & H_p(\text{Alpha}(s)). \end{array}$$
- Prove that the vertical maps are isomorphisms and that the diagram commutes.
- Show that the persistence diagrams of f are the same as those of the sequence of alpha complexes.

References

- [1] E. A. ABBOT. *Flatland*. Dover, New York, 1952.
- [2] C. C. ADAMS. *The Knot Book. An Elementary Introduction to the Mathematical Theory of Knots*. Amer. Math. Soc., Providence, Rhode Island, 2004.
- [3] P. K. AGARWAL, H. EDELSBRUNNER, J. L. HARER AND Y. WANG. Extreme elevation on a 2-manifold. *Discrete Comput. Geom.* **36** (2006), 553–572.
- [4] P. K. AGARWAL, H. EDELSBRUNNER AND Y. WANG. Computing the writhing number of a polygonal knot. *Discrete Comput. Geom.* **32** (2004), 37–53.
- [5] P. K. AGARWAL, A. EFRAT AND M. SHARIR. Vertical decomposition of shallow levels in 3-dimensional arrangements and its applications. *SIAM J. Comput.* **29** (2000), 912–953.
- [6] A. V. AHO, J. E. HOPCROFT AND J. D. ULLMAN. *The Design and Analysis of Computer Algorithms*. Addison-Wesley, Reading, Massachusetts, 1973.
- [7] R. K. AHUJA, T. L. MAGNANTI AND J. B. ORLIN. *Network Flows: Theory, Algorithms, and Applications*. Prentice Hall, Englewood Cliffs, New Jersey, 1993.
- [8] J. W. ALEXANDER. A proof of the invariance of certain constants of analysis situ. *Trans. Amer. Math. Soc.* **16** (1915), 148–154.
- [9] P. S. ALEXANDROV. Über den allgemeinen Dimensionsbegriff und seine Beziehungen zur elementaren geometrischen Anschauung. *Math. Ann.* **98** (1928), 617–635.
- [10] V. I. ARNOLD. *Ordinary Differential Equations*. Translated from Russian, MIT Press, Cambridge, Massachusetts, 1973.
- [11] F. AURENHAMMER. Voronoi diagrams — a study of a fundamental geometric data structure. *ACM Comput. Surveys* **23** (1991), 345–405.
- [12] T. F. BANCHOFF. Critical points and curvature for embedded polyhedra. *J. Differential Geometry* **1** (1967), 245–256.
- [13] T. F. BANCHOFF. Triple points and surgery of immersed surfaces. *Proc. Amer. Math. Soc.* **46** (1974), 403–413.
- [14] A. BANYAGA AND D. HURTUBIS. *Lectures on Morse Homology*. Kluwer, Dordrecht, The Netherlands, 2004.
- [15] W. R. BAUER, F. H. C. CRICK AND J. H. WHITE. Supercoiled DNA. *Scientific American* **243** (1980), 118–133.
- [16] B. BAUMGART. A polyhedron representation for computer vision. In “Proc. Natl. Comput. Conf., 1975”, 589–596.
- [17] P. BENDICH, D. COHEN-STEINER, H. EDELSBRUNNER, J. L. HARER AND D. MOROZOV. Inferring local homology from sampled stratified spaces. In “Proc. 48th Ann. Sympos. Found. Comput. Sci., 2007”, 536–546.
- [18] S. BEUCHER. Watersheds of functions and picture segmentation. In “Proc. IEEE Intl. Conf. Acoustic, Speech, Signal Process, 1982”, 1928–1931.
- [19] K. BORSUK. On the imbedding of systems of compacta in simplicial complexes. *Fund. Math.* **35** (1948), 217–234.

- [20] H. R. BRAHANA. Systems of circuits on two-dimensional manifolds. *Ann. Math.* **23** (1922), 144–168.
- [21] E. BRISSON. Representing geometric structures in d dimensions: topology and order. *Discrete Comput. Geom.* **9** (1993), 387–426.
- [22] L. E. J. BROUWER. Über eineindeutige, stetige Transformationen von Flächen in sich. *Math. Ann.* **69** (1910), 176–180.
- [23] L. E. J. BROUWER. Über Abbildungen von Mannigfaltigkeiten. *Math. Ann.* **71** (1912), 97–115.
- [24] K. S. BROWN. *Cohomology of Groups*. Springer-Verlag, New York, New York, 1994.
- [25] G. CĂLUGĂREANU. Sur les classes d'isotopie des noeuds tridimensionnels et leurs invariants. *Czech. Math. J.* **11** (1961), 588–625.
- [26] W. A. CANNON. A tentative classification of root systems. *Ecology* **30** (1947), 452–458.
- [27] H. CARR, J. SNOEYINK AND U. AXEN. Computing contour trees in all dimensions. *Comput. Geom. Theory Appl.* **24** (2002), 75–94.
- [28] J. S. CARTER. *How Surfaces Intersect in Space. An Introduction to Topology*. Second edition, World Scientific, Singapore, 1995.
- [29] F. CAZALS, F. CHAZAL AND T. LEWINER. Molecular shape analysis based upon the Morse-Smale complex and the Connolly function. In “Proc. 19th Ann. Sympos. Comput. Geom., 2003”, 237–246.
- [30] J. CERF. La stratification naturelle des espaces de fonctions différentiables réelles et le théorème de la pseudo-isotopie. *Inst. Hautes Etudes Sci. Publ. Math.* **39** (1970), 5–173.
- [31] F. CHAZAL, D. COHEN-STEINER, M. GLISSE, L. J. GUIBAS AND S. Y. OUDOT. Proximity of persistence modules and their diagrams. In “Proc. 25th Ann. Sympos. Comput. Geom., 2009”, 237–246.
- [32] B. CHAZELLE. Triangulating a simple polygon in linear time. *Discrete Comput. Geom.* **6** (1991), 485–524.
- [33] D. COHEN-STEINER AND H. EDELSBRUNNER. Inequalities for the curvature of curves and surfaces. *Found. Comput. Math.* **7** (2007), 391–404.
- [34] D. COHEN-STEINER, H. EDELSBRUNNER AND J. L. HARER. Stability of persistence diagrams. *Discrete Comput. Geom.* **37** (2007), 103–120.
- [35] D. COHEN-STEINER, H. EDELSBRUNNER AND J. L. HARER. Extended persistence using Poincaré and Lefschetz duality. *Found. Comput. Math.* **9** (2009), 79–103. Erratum. *Found. Comput. Math.* **9** (2009), 133–134.
- [36] D. COHEN-STEINER, H. EDELSBRUNNER, J. L. HARER AND Y. MILEYKO. Lipschitz functions have L_p -stable persistence. *Found. Comput. Math.*, to appear.
- [37] D. COHEN-STEINER, H. EDELSBRUNNER, J. L. HARER AND D. MOROZOV. Persistent homology for kernels, images, and cokernels. In “Proc. 20th Ann. ACM-SIAM Sympos. Discrete Alg., 2009”, 1011–1020.
- [38] D. COHEN-STEINER, H. EDELSBRUNNER AND D. MOROZOV. Vines and vineyards by updating persistence in linear time. In “Proc. 22nd Ann. Sympos. Comput. Geom., 2006”, 119–126.
- [39] K. COLE-MCLAUGHLIN, H. EDELSBRUNNER, J. L. HARER, V. NATARAJAN AND V. PASCUCCI. Loops in Reeb graphs of 2-manifolds. *Discrete Comput. Geom.* **32** (2004), 231–244.
- [40] M. L. CONNOLLY. Shape complementarity at the hemo-globin albl subunit interface. *Biopolymers* **25** (1986), 1229–1247.
- [41] T. H. CORMEN, C. E. LEISERSON, R. L. RIVEST AND C. STEIN. *Introduction to Algorithms*. Second edition, McGraw Hill, Boston, Massachusetts, 2001.
- [42] W. J. CROFT. *Under the Microscope. A Brief History of Microscopy*. World Scientific, Singapore, Singapore, 2006.
- [43] M. DEHN AND P. HEGGARD. Analysis situ. *Enzykl. Math. Wiss.* **III AB** **3** (1907), 153–220.

- [44] B. DELAUNAY. Sur la sphère vide. *Izv. Akad. Nauk SSSR, Otdelenie Matematicheskii i Estestvennyka Nauk* **7** (1934), 793–800.
- [45] C. J. A. DELFINADO AND H. EDELSBRUNNER. An incremental algorithm for Betti numbers of simplicial complexes on the 3-sphere. *Comput. Aided Geom. Design* **12** (1995), 771–784.
- [46] M.-L. DEQUÈANT, S. AHNERT, H. EDELSBRUNNER, T. M. A. FINK, E. F. GLYNN, G. HATTEM, A. KUDLICKI, Y. MILEYKO, J. MORTON, A. R. MUSHEGIAN, L. PACHTER, M. ROWICKA, A. SHIU, B. STURMFELS AND O. POURQUIÉ. Comparison of pattern detection methods in microarray time series of the segmentation clock. *PLoS ONE* **3** (2008), e2856, doi:10.1371/journal.pone.0002856.
- [47] M.-L. DEQUÈANT, E. F. GLYNN, K. GAUDENZ, M. WAHL, J. CHEN, A. R. MUSHEGIAN AND O. POURQUIÉ. A complex oscillating network of signaling genes underlies the mouse segmentation clock. *Science* **314** (2006), 1595–1598.
- [48] T. DEY, H. EDELSBRUNNER, S. GUHA AND D. V. NEKHAYEV. Topology preserving edge contraction. *Publ. Inst. Math. (Beograd) (N.S.)* **66** (1999), 23–45.
- [49] E. W. DIJKSTRA. A note on two problems in connexion with graphs. *Numerische Mathematik* **1** (1959), 269–271.
- [50] E. A. DINIC. Algorithm for solution of a problem of maximum flow in a network with power estimation. *Soviet Math. Doklady* **11** (1970), 1277–1280.
- [51] D. P. DOBKIN AND M. J. LASZLO. Primitives for the manipulation of three-dimensional subdivisions. *Algorithmica* **4** (1989), 3–32.
- [52] J. DRENTH. *Principles of Protein X-Ray Crystallography*. Springer-Verlag, New York, New York, 1999.
- [53] H. EDELSBRUNNER. The union of balls and its dual shape. *Discrete Comput. Geom.* **13** (1995), 415–440.
- [54] H. EDELSBRUNNER. *Geometry and Topology for Mesh Generation*. Cambridge Univ. Press, Cambridge, England, 2001.
- [55] H. EDELSBRUNNER. Surface tiling with differential topology. In “Proc. 3rd Eurographics Sympos. Geom. Process., 2005”, 9–11.
- [56] H. EDELSBRUNNER AND L. J. GUIBAS. Topologically sweeping an arrangement. *J. Comput. System Sci.* **38** (1989), 165–194. Corrigendum. *J. Comput. System Sci.* **42** (1991), 249–251.
- [57] H. EDELSBRUNNER AND J. L. HARER. Persistent homology — a survey. *Surveys on Discrete and Computational Geometry. Twenty Years Later*, J. E. Goodman, J. Pach and R. Pollack (eds.), Contemporary Mathematics **453**, 257–282, Amer. Math. Soc., Providence, Rhode Island, 2008.
- [58] H. EDELSBRUNNER, J. L. HARER AND A. J. ZOMORODIAN. Hierarchical Morse-Smale complexes for piecewise linear 2-manifolds. *Discrete Comput. Geom.* **30** (2003), 87–107.
- [59] H. EDELSBRUNNER, D. G. KIRKPATRICK AND R. SEIDEL. On the shape of a set of points in the plane. *IEEE Trans. Inform. Theory* **IT-29** (1983), 551–559.
- [60] H. EDELSBRUNNER, D. LETSCHER AND A. J. ZOMORODIAN. Topological persistence and simplification. *Discrete Comput. Geom.* **28** (2002), 511–533.
- [61] H. EDELSBRUNNER, D. MOROZOV AND V. PASCUCCI. Persistence-sensitive simplification of functions on 2-manifolds. In “Proc. 22nd Ann. Sympos. Comput. Geom., 2006”, 127–134.
- [62] H. EDELSBRUNNER AND E. P. MÜCKE. Simulation of Simplicity: a technique to cope with degenerate cases in geometric algorithms. *ACM Trans. Graphics* **9** (1990), 86–104.
- [63] H. EDELSBRUNNER AND E. P. MÜCKE. Three-dimensional alpha shapes. *ACM Trans. Graphics* **13** (1994), 43–72.
- [64] R. D. EDWARDS. Approximating certain cell-like maps by homeomorphisms. *Notices Amer. Math. Soc.* **24** (1977), A647.
- [65] A. EFRAT, A. ITAI AND M. J. KATZ. Geometry helps in bottleneck matching and related problems. *Algorithmica* **31** (2001), 1–28.

- [66] S. EILENBERG AND N. STEENROD. *Foundations of Algebraic Topology*. Princeton Univ. Press, Princeton, New Jersey, 1952.
- [67] I. FÁRY. Sur certaines inégalités géométriques. *Acta Sci. Math. Szeged* **12** (1950), 117–124.
- [68] M. S. FLOATER. One-to-one piecewise linear mappings over triangulations. *Math. Comput.* **72** (2003), 685–696.
- [69] A. FLOER. Witten’s complex and infinite dimensional Morse theory. *J. Diff. Geom.* **30** (1989), 207–221.
- [70] A. FLORES. Über n -dimensionale Komplexe die in \mathbb{R}_{2n+1} selbstverschlungen sind. *Ergeb. Math. Koll.* **6** (1933/34), 4–7.
- [71] A. R. FOREST. Computational geometry in practice. In *Fundamental Algorithms for Computer Graphics*, E. A. Earnshaw (ed.), Springer-Verlag, Berlin, Germany, 1985, 707–724.
- [72] D. A. FORSYTH AND J. PONCE. *Computer Vision: a Modern Approach*. Prentice Hall, Englewood Cliffs, New Jersey, 2003.
- [73] P. FROSINI AND C. LANDI. Size theory as a topological tool for computer vision. *Pattern Recognition and Image Analysis* **9** (1999), 596–603.
- [74] F. B. FULLER. The writhing number of a space curve. *Proc. Natl. Acad. Sci. USA* **68** (1971), 815–819.
- [75] M. GARLAND AND P. S. HECKBERT. Surface simplification using quadric error metrics. *Computer Graphics*, Proc. SIGGRAPH, 1997, 209–216.
- [76] I. M. GELFAND, M. M. KAPRANOV AND A. V. ZELEVINSKY. *Discriminants, Resultants and Multidimensional Determinants*. Birkhäuser, Boston, Massachusetts, 1994.
- [77] L. GEORGIADIS, R. E. TARJAN AND R. F. WERNECK. Design of data structures for mergeable trees. In “Proc. 17th Ann. ACM-SIAM Sympos. Discrete Alg., 2006”, 394–403.
- [78] P. J. GIBLIN. *Graphs, Surfaces and Homology*. Chapman and Hall, London, England, 1981.
- [79] M. GROMOV. Hyperbolic groups. In *Essays in Group Theory*, 75–263, Math. Sci. Res. Inst. Publ. **8**, Springer-Verlag, New York, New York, 1987.
- [80] L. J. GUIBAS AND J. STOLFI. Primitives for the manipulation of general subdivisions and the computation of Voronoi diagrams. *ACM Trans. Graphics* **4** (1985), 74–123.
- [81] E. M. HAACKE, R. W. BROWN, M. R. THOMPSON AND R. VENKATESAN. *Magnetic Resonance Imaging: Physical Principles and Sequence Design*. Wiley, New York, New York, 1999.
- [82] A. HATCHER. *Algebraic Topology*. Cambridge Univ. Press, Cambridge, England, 2002.
- [83] E. HELLY. Über Mengen konvexer Körper mit gemeinschaftlichen Punkten. *Jahresber. Deutsch. Math.-Verein.* **32** (1923), 175–176.
- [84] E. HELLY. Über Systeme von abgeschlossenen Mengen mit gemeinschaftlichen Punkten. *Monatsh. Math. Physik* **37** (1930), 281–302.
- [85] J. E. HOPCROFT AND R. M. KARP. An $n^{5/2}$ algorithm for maximum matchings in bipartite graphs. *SIAM J. Comput.* **2** (1973), 225–231.
- [86] B. K. P. HORN. *Robot Vision*. MIT Press, Cambridge, Massachusetts, 1986.
- [87] J. JANIN, K. HENRICK, J. MOULT, L. T. EYCK, M. J. STERNBERG, S. VAJDA, I. VAKSER AND S. J. WODAK. CAPRI: a critical assessment of predicted interactions. *Proteins* **52** (2003), 2–9.
- [88] R. KANNAN AND A. BACHEM. Polynomial algorithms for computing the Smith and Hermite normal forms of an integer matrix. *SIAM J. Comput.* **8** (1979), 499–507.
- [89] L. V. KANTOROVICH. On the translocation of masses. *C. R. (Dokl.) Acad. Sci. USSR* **37** (1942), 199–226.
- [90] J. L. KELLEY AND E. PITCHER. Exact homomorphism sequences in homology theory. *Ann. of Math.* **48** (1947), 682–709.
- [91] D. P. KIEHART, C. G. GALBRAITH, K. A. EDWARDS, W. L. RICKOLL AND R. A. MONTAGUE. Multiple forces contribute to cell sheet morphogenesis for dorsal closure in *Drosophila*. *J. Cell Biology* **149** (2000), 471–490.

- [92] J. KLEINBERG AND E. TARDOS. *Algorithm Design*. Pearson Education, Boston, Massachusetts, 2006.
- [93] R. KLETTE AND A. ROSENFIELD. *Digital Geometry. Geometric Methods for Digital Picture Analysis*. Morgan Kaufmann, San Francisco, California, 2004.
- [94] D. E. KNUTH. *Sorting and Searching. The Art of Computer Programming, Vol. 3*. Addison-Wesley, Reading, Massachusetts, 1973.
- [95] P. KOEBE. Kontaktprobleme der konformen Abbildung. *Ber. Sächs. Akad. Wiss. Leipzig, Math.-Phys. Kl.* **88** (1936), 141–164.
- [96] H. W. KUHN. The Hungarian Method for the assignment problem. *Naval Research Logistics Quarterly* **2** (1955), 83–97.
- [97] I. KUPKA. Contribution à la théorie des champs génériques. *Contributions to Differential Equations* **2** (1963), 457–484.
- [98] K. N. KUTULAKOS AND S. M. SEITZ. A theory of shape by space carving. *Internat. J. Comput. Vision* **38** (2000), 199–218.
- [99] I. LAKATOS. *Proofs and Refutations: the Logic of Mathematical Discovery*. Cambridge Univ. Press, Cambridge, England, 1976.
- [100] A. LAURENTINI. The visual hull concept for silhouette-based image understanding. *IEEE Trans. Pattern Anal. Mach. Intell.* **PAMI-16** (1994), 150–162.
- [101] S. LEFSCHETZ. Manifolds with a boundary and their transformations. *Trans. Amer. Math. Soc.* **29** (1927), 429–462.
- [102] J. LERAY. Sur la forme des espaces topologiques et sur les points fixes des représentations. *J. Math. Pures Appl.* **24** (1945), 95–167.
- [103] W. E. LORENSEN AND H. E. CLINE. Marching cubes: a high resolution 3D surface construction algorithm. *Comput. Graphics* **21**, Proc. SIGGRAPH, 1987, 163–169.
- [104] A. A. MARKOV. Insolubility of the problem of homeomorphy. In *Proc. Int. Congr. Math.*, 1958, 14–21.
- [105] Y. MATSUMOTO. *An Introduction to Morse Theory*. Translated from Japanese by K. Hudson and M. Saito, Amer. Math. Soc., Providence, Rhode Island, 2002.
- [106] C. R. F. MAUNDER. *Algebraic Topology*. Cambridge Univ. Press, Cambridge, England, 1980.
- [107] J. P. MAY. *A Concise Course in Algebraic Topology*. Univ. Chicago Press, Chicago, Illinois, 1999.
- [108] W. MAYER. Über abstrakte Topologie. *Monatschr. Math. Phys.* **36** (1929), 1–42 and 219–258.
- [109] J. McCLEARY. *A User's Guide to Spectral Sequences*. Second edition, Cambridge Univ. Press, Cambridge, England, 2001.
- [110] J. MILNOR. Two complexes which are homeomorphic but combinatorially distinct. *Ann. of Math.* **74** (1961), 575–590.
- [111] J. MILNOR. *Morse Theory*. Princeton Univ. Press, Princeton, New Jersey, 1963.
- [112] G. MONGE. Mémoire sur la théorie des déblais et des remblais. In *Histoire de l'Académie Royale des Sciences de Paris* (1781), 666–704.
- [113] M. MORSE. *The Calculus of Variations in the Large*. Amer. Math. Soc., New York, New York, 1934.
- [114] M. MORSE. Topologically non-degenerate functions on a compact n -manifold. *J. Analyse Math.* **7** (1959), 189–208.
- [115] J. R. MUNKRES. *Topology. A First Course*. Prentice-Hall, Englewood Cliffs, New Jersey, 1975.
- [116] J. R. MUNKRES. *Elements of Algebraic Topology*. Perseus, Cambridge, Massachusetts, 1984.
- [117] T. NISHIZEKI AND N. CHIBA. *Planar Graphs: Theory and Algorithms*. North-Holland, Amsterdam, The Netherlands, 1988.

- [118] I. PALMEIRIM, D. HENRIQUE, D. ISH-HOROWICZ AND O. POURQUIÉ. Avian hairy gene expression identifies a molecular clock linked to vertebrate segmentation and somitogenesis. *Cell* **91** (1997), 639–648.
- [119] W. F. POHL. The self-linking number of a closed space curve. *J. Math. Mech.* **17** (1968), 975–985.
- [120] H. POINCARÉ. Analysis situs. *J. Ecole Polytechn.* **1** (1895), 1–121.
- [121] H. POINCARÉ. Complément à l'analysis situs. *Rend. Circ. Mat. Palermo* **13** (1899), 285–343.
- [122] H. POINCARÉ. Cinquième complément à l'analysis situs. *Rend. Circ. Mat. Palermo* **18** (1904), 45–110.
- [123] O. POURQUIÉ. The segmentation clock: converting embryonic time into spatial pattern. *Science* **301** (2003), 328–330.
- [124] F. P. PREPARATA AND M. I. SHAMOS. *Computational Geometry: an Introduction*. Springer-Verlag, New York, 1985.
- [125] A. A. RANICKI (ED.). *The Hauptvermutung Book*. Kluwer, Dordrecht, The Netherlands, 1996.
- [126] G. REEB. Sur les points singuliers d'une forme de Pfaff complètement intégrable ou d'une fonction numérique. *Comptes Rendus de L'Académie ses Séances, Paris* **222** (1946), 847–849.
- [127] V. ROBINS. Toward computing homology from finite approximations. *Topology Proceedings* **24** (1999), 503–532.
- [128] J. ROERDINK AND A. MEIJSTER. The watershed transform: definitions, algorithms, and parallelization strategies. *Fundamenta Informaticae* **41** (2000), 187–228.
- [129] L. SANTALÓ. *Integral Geometry and Geometric Probability*. Addison-Wesley, 1976, reprinted by Cambridge Univ. Press, Cambridge, England, 2004.
- [130] T. SCHLICK. *Molecular Modeling and Simulation. An Interdisciplinary Guide*. Springer-Verlag, New York, New York, 2002.
- [131] J. SIJBERS, P. SCHEUNDERS, M. VERHOYE, A. VAN DER LINDEN, D. VAN DYCK AND E. RAMAN. Watershed-based segmentation of 3D MR data for volume quantization. *Magn. Reson. Imag.* **15** (1997), 679–688.
- [132] D. D. SLEATOR AND R. E. TARJAN. Self-adjusting binary search trees. *J. Assoc. Comput. Mach.* **32** (1985), 652–686.
- [133] S. SMALE. Stable manifolds for differential equations and diffeomorphisms. *Ann. Scuola Norm. Sup. Pisa* **17** (1963), 97–116.
- [134] H. J. SMITH. On systems of indeterminate equations and congruences. *Philos. Trans.* **151** (1861), 293–326.
- [135] M. SONKA, V. HLAVAC AND R. BOYLE. *Image Processing, Analysis and Machine Vision*. Second edition, PWS Publishing, Pacific Grove, California, 1999.
- [136] E. H. SPANIER. *Algebraic Topology*. Springer-Verlag, New York, New York, 1966.
- [137] E. STEINITZ. *Polyeder und Raumeinteilung*. In *Enzykl. Math. Wiss.* **III AB** **12** (1922), 1–139.
- [138] A. STORJOHANN. Near optimal algorithm for computing Smith normal forms of integer matrices. In “Proc. Internat. Sympos. Symbol. Algebraic Comput., 1997”, 267–274.
- [139] G. STRANG. *Introduction to Linear Algebra*. Wellesley-Cambridge Press, Wellesley, Massachusetts, 1993.
- [140] R. E. TARJAN. *Data Structures and Network Algorithms*. SIAM, Philadelphia, Pennsylvania, 1983.
- [141] W. T. TUTTE. How to draw a graph. *Proc. London Math. Soc.* **13** (1963), 743–768.
- [142] W. T. TUTTE. *Graph Theory*. Addison-Wesley, Reading, Massachusetts, 1984.

- [143] E. R. VAN KAMPEN. Komplexe in euklidischen Räumen. *Abh. Math. Sem. Univ. Hamburg* **9** (1933), 72–78.
- [144] M. VAN KREVELD, R. VAN OOSTRUM, C. L. BAJAJ, V. PASCUCCI AND D. R. SCHIKORE. Contour trees and small seed sets for isosurface traversal. In “Proc. 13th Ann. Sympos. Comput. Geom., 1997”, 212–220.
- [145] O. VEBLEN. Theory on plane curves in non-metrical analysis situs. *Trans. Amer. Math. Soc.* **6** (1905), 83–98.
- [146] L. VIETORIS. Über den höheren Zusammenhang kompakter Räume und eine Klasse von zusammenhangstreuen Abbildungen. *Math. Ann.* **97** (1927), 454–472.
- [147] L. VIETORIS. Über die Homologiegruppen der Vereinigung zweier Komplexe. *Monatsh. Math.* **37** (1930), 159–162.
- [148] C. VILLANI. *Topics in Optimal Transportation*. Amer. Math. Soc., Providence, Rhode Island, 2003.
- [149] G. VORONOI. Nouvelles applications des paramètres continus à la théorie des formes quadratiques. Premier Mémoire: Sur quelques propriétés des formes quadratiques positives parfaites. *J. Reine Angew. Math.* **133** (1907), 97–178.
- [150] G. VORONOI. Nouvelles applications des paramètres continus à la théorie des formes quadratiques. Deuxième Mémoire: Recherches sur les paralléloèdres primitifs. *J. Reine Angew. Math.* **134** (1908), 198–287.
- [151] C. T. C. WALL. *A Geometric Introduction to Topology*. Addison-Wesley, Reading, Massachusetts, 1972.
- [152] Y. WANG, P. K. AGARWAL, P. BROWN, H. EDELSBRUNNER AND J. RUDOLPH. Coarse and reliable geometric alignment for protein docking. In “Proc. Pacific Sympos. Biocomput., 2005”, 65–75.
- [153] L. N. WASSERSTEIN. Markov processes over denumerable products of spaces describing large systems of automata. *Problems of Information Transmission* **5** (1969), 47–52.
- [154] J. D. WATSON AND F. H. C. CRICK. Molecular structure of nucleic acid. A structure for deoxyribose nucleic acid. *Nature* **171** (1953), 737–738.
- [155] E. WELZL. Smallest enclosing disks (balls and ellipsoids). In *New Results and New Trends in Computer Science*, H. A. Maurer (ed.), Springer-Verlag, Lecture Notes in Computer Science **555** (1991), 359–370.
- [156] J. H. WHITE. Self-linking and the Gauss integral in higher dimensions. *Amer. J. Math.* **XCI** (1969), 693–728.
- [157] H. WHITNEY. The self-intersections of a smooth n -manifold in $2n$ -space. *Annals of Math.* **45** (1944), 220–246.
- [158] H. WHITNEY. The singularities of a smooth n -manifold in $(2n - 1)$ -space. *Annals of Math.* **45** (1944), 247–293.
- [159] K. WUTHRICH. *NMR of Proteins and Nucleic Acids*. Wiley-Interscience, New York, New York, 1986.
- [160] A. J. ZOMORODIAN. *Topology for Computing*. Cambridge Univ. Press, Cambridge, England, 2005.
- [161] A. J. ZOMORODIAN AND G. CARLSSON. Computing persistent homology. *Discrete Comput. Geom.* **33** (2005), 249–274.

Index

- absolute homology group, 90
- abstract simplicial complex, 53
- Ackermann function, 8, 123, 145, 159
- affine
 - combination, 51
 - hull, 51
 - independence, 51
- Alexander
 - duality, 118
 - Duality Theorem, 120
- alpha
 - complex, 69, 223
 - weighted, 69
 - shape, 69, 223
- anti-cancellation, 208
- ascending manifold, 131
- augmentation map, 83
- augmenting path, 192
- balanced search tree, 144
- barycenter, 54
- barycentric
 - coordinate, 54, 135
 - subdivision, 54, 74, 119
- basis of a topology, 5
- Betti number, 81, 85
 - persistent, 151, 182
 - reduced, 83
- bipartite graph, 191
- birth, 151
- bisector, 65
- block, 168
 - chain, 111
- Block Complex Lemma, 111
- body centered cubic lattice, 146
- Borromean rings, 24
- bottleneck distance, 180
- boundary, 80
 - group, 80
 - relative, 90
- homomorphism, 80
 - map, 80, 95, 169
 - matrix, 86, 167
 - of a manifold, 28
 - of a simplex, 52
- branch, 220
- branch point, 42
- Breadth-first Search, 35, 140, 193
- Brouwer's Fixed Point Theorem, 84
- Călugăreanu-White Formula, 17
- cancellation, 208
- Cauchy-Crofton Formula, 187
- Čech complex, 60
- chain, 79
 - complex, 80, 95, 169
 - group, 80
 - relative, 90
 - map, 95
- Classification Theorem for 2-manifolds, 29
- closed
 - curve, 9, 186
 - simple, 9, 208
 - polygon, 10, 210
 - set, 6
 - star, 53
- coboundary, 105
 - group, 105
 - map, 105
 - matrix, 107
- cochain, 104
 - group, 105
- cocycle, 105
 - group, 105
- coface, 52
- coherent triangulation, 68
- cohomology group, 105
 - reduced, 105
- cokernel, 93
- collapse, 72
- collapsible, 72, 75
- collision, 157
- coloring, 14, 24, 48
- combination
 - affine, 51
 - convex, 20, 51
- commutative square, 97
- compact, 27
- compatible ordering, 152, 167
- complementary subcomplexes, 119
- complex
 - abstract simplicial, 53
 - simplicial, 52
- component, 4, 158

- connected, 4
 - sum, 28
- connecting homomorphism, 94, 96
- continuous, 5
- contour, 141
 - tree, 140
- contractible, 59
- contraction of an edge, 42
- convex
 - combination, 20, 51
 - hull, 20, 51
 - set system, 57
- coordinate chart, 37
- coset, 81
- cost function, 191
- critical
 - event, 73
 - point, 127
 - value, 127
 - homological, 183
 - vertex, 136
- cross-cap, 29
- curvature, 186
 - Gaussian, 197
 - mean, 197, 212
 - total, 186
- curve, 9, 186, 208
- cycle, 80
 - group, 80
 - relative, 90
- cyclic list, 143
- cylinder, 28
- death, 151
- decidability, 32
- deduction map, 195
- deformation
 - retract, 58
 - retraction, 58
- Delaunay
 - complex, 67, 69
 - weighted, 68
 - triangulation, 67
- density data, 140
- Depth-first Search, 35, 140, 193
- descending manifold, 131
- destination, 130
- diameter, 55
- diffeomorphism, 126
- Dijkstra's Algorithm, 195
- dimension
 - of a complex, 52
 - of a simplex, 52
- direct sum, 93
- directed graph, 195
- directional writhing number, 16
- Dirichlet tessellation, 68
- disjoint set system, 6
- disk, 28
- distance
 - bottleneck, 180
 - Fréchet, 188
 - power, 65
 - signed, 45
 - squared, 45
 - weighted, 65
- Wasserstein, 183
- distance function, 223
- doubling of a manifold, 32
- drawing of a graph, 18
- dual
 - block, 110
 - decomposition, 110, 119
 - homomorphism, 104
- duality, 163
 - Alexander, 118
 - Lefschetz, 117, 162, 166
 - Poincaré, 109, 114, 166
- dunce cap, 102
- edge contraction, 42
- Elder Rule, 150
- elementary collapse, 72
- elevation
 - function, 162, 209
 - maximum
 - one-legged, 210
 - two-legged, 210
- embedding, 13, 38
 - of a graph, 18
 - straight-line, 18
- equivalence of knots, 13
- essential, 162
- Euler
 - characteristic, 31, 85, 133
 - Characteristic of 2-manifolds, 31
 - Poincaré Theorem, 86, 137
 - Relation for Planar Graphs, 18
- event
 - critical, 73
 - regular, 73
- exact, 93, 95
 - sequence
 - of a pair, 94
 - of a triple, 102
 - of chain complexes, 95
- Exact Sequence of a Pair Theorem, 94
- Excision Theorem, 92
- Existence and Uniqueness Theorems, 131
- extended
 - persistence, 161, 208
 - diagram, 163
 - real plane, 152
- Extended Persistence Algorithm, 213
- exterior, 119
- face
 - of a planar graph, 18

- of a simplex, 52
- Fáry Theorem, 188
- field, 156
- figure-eight knot, 13
- filtration, 70, 151, 166
 - lower star, 135, 159
- First
 - Plane Lemma, 67
 - Sphere Lemma, 66
- fixed point, 84, 102
- flag, 74
- Floer homology, 134
- fork, 220
- formal sum, 79
- Fréchet distance, 188
- full subcomplex, 52
- function
 - elevation, 162, 209
 - height, 125, 161
 - monotonic, 175
 - Morse, 208
 - PL, 135, 160
 - smooth, 126
- Fundamental Lemma
 - of Homology, 81
 - of Persistent Homology, 152
- fundamental quadric, 46
- Gaussian
 - curvature, 197
 - elimination, 88
- gene expression, 200
- general position, 67
- Generalized Fáry Theorem, 188
- generic
 - curve, 210
 - PL function, 135
- genotype, 218
- genus, 31, 36, 142
- geometric realization, 53
- Geometric Realization Theorem, 53
- gradient, 38, 129
- graph
 - abstract, 3
 - bipartite, 191
 - complete, 191
 - coloring, 24, 48
 - complete, 3
 - directed, 195
 - homeomorphism, 20
 - matching, 191
 - planar, 18
 - maximally connected, 19
 - Reeb, 141
 - simple, 3
 - weighted, 196
- group
 - of boundaries, 80
 - of chains, 80
- of coboundaries, 105
- of cochains, 105
- of cocycles, 105
- of cycles, 80
- of diffeomorphisms, 129
- of homomorphisms, 104
- Hasse diagram, 71
- Hauptvermutung, 57
- height function, 125, 161
- Helly's Theorem, 57
- Hessian, 127
- homeomorphism, 9
- homological critical value, 183
- homologous, 82
- homology
 - class, 81
 - group, 81
 - absolute, 90
 - persistent, 151, 182
 - reduced, 83, 176
 - relative, 90, 162, 170
 - local, 220
 - homomorphism, 80
 - homotopic, 58
 - homotopy, 58
 - equivalence, 58
 - inverse, 58
 - straight-line, 181
 - type, 58
 - Hopf link, 15
 - hull
 - affine, 51
 - convex, 20, 51
 - image, 93
 - analysis, 213
 - segmentation, 214
 - immersion, 38
 - Incremental Betti Number Algorithm, 121
 - index
 - of a critical point, 128
 - of a PL critical vertex, 136
 - index persistence, 152
 - integral
 - geometry, 187
 - line, 130
 - interchange, 208
 - interior of a simplex, 52
 - intersection
 - number, 114
 - pairing, 115
 - inversion, 63
 - Inversion Lemma, 63
 - irreducible triangulation, 48
 - iso-surface, 140
 - isomorphism
 - between complexes, 53, 54, 74
 - between homology groups, 92

- Iteration Bound, 194
- Jacobian, 38
- join, 109
- Jordan Curve Theorem, 9
- Jung's Theorem, 60
- kernel, 93
- Klein bottle, 29, 41, 102, 116
- knot, 13
- Kuratowski Theorem, 20
- Lefschetz
 - duality, 117, 162, 166
 - Duality Theorem, 117, 118
- length, 187
- level set, 125, 141
- lifting, 65
- limit term, 171
- line arrangement, 180
- linear
 - array, 6, 34, 156
 - equation, 23
- link
 - of a simplex, 53
 - of an edge, 44
 - of knots, 15
- Link Condition Lemma, 44
- linked list, 156
- linking number, 15
- Lipschitz, 184, 201
- list, 143
- local homology, 220
- long exact sequence, 94
- longitudinal curve, 116
- loop in a Reeb graph, 141
- Loop Lemma for Manifolds, 142
- lower
 - link, 136
 - star, 135, 165
 - filtration, 135, 159
- lowest 1, 153
- Möbius strip, 28, 41
- magnetic resonance imaging, 214
- manifold, 141, 160
 - ascending, 131
 - descending, 131
 - stable, 131
 - unstable, 131
 - with boundary, 28
 - without boundary, 27
- Marching Cube Algorithm, 145
- matching, 191
 - maximum, 191
 - minimum cost, 191
 - perfect, 191
- matrix, 176
 - boundary, 86, 167
- decomposition, 176
- reduction, 88
- sparse, 156
- maximum, 128
 - matching, 191
 - principle, 21
- Maximum Matching Algorithm, 196
- Mayer-Vietoris
 - sequence, 98, 138
 - Sequence Theorem, 98
- mean curvature, 197, 212
 - total, 212
- measure, 203
 - theory, 175
- merge tree, 149
- meridian curve, 116
- mesh, 55, 184
- mesh generation, 68
- metric, 181
- microarray, 200
- microscopy, 213
- miniball, 60
- minimum, 128
 - cost matching, 191
- Minimum Cost Matching Algorithm, 196
- molecular
 - skin, 208
 - surface, 208
- monkey saddle, 137
- monotonic function, 150, 175
- Morse
 - function, 128, 208
 - topological, 74
 - Inequalities, 133
 - Lemma, 127
 - Smale
 - complex, 133
 - function, 132
 - Witten complex, 134
- multiplicity, 152, 183
- multiset, 152
- negative simplex, 121, 154
- nerve, 59, 67, 69
- Nerve Theorem, 59
- non-degenerate critical point, 127
- non-orientable, 29
- nuclear magnetic resonance, 207
- open
 - cover, 27
 - set, 4
- optimal transportation, 185
- order
 - complex, 74
 - of a group, 82
- ordered triangle, 33
- ordinary persistence diagram, 163
- orientable, 29

- orientation, 33, 39
 - preserving, 28
 - reversing, 28
- origin, 130
- output-sensitive, 158
- pair of spaces, 90
- pairing, 115, 153
 - intersection, 115
 - perfect, 115
- Pairing Lemma, 154
- Parity Algorithm, 10
- path, 5, 9
 - augmenting, 192
 - compression, 8
 - connected, 5
 - decomposition, 150
 - shortest, 196
- perfect
 - matching, 191
 - pairing, 115
- periodicity, 202
- Periodicity Measure Lemma, 203
- Persistence
 - Algorithm, 167, 216
 - Duality Theorem, 164
 - Equivalence Theorem, 159
 - Symmetry Theorem, 164
- persistence, 152
 - diagram, 152, 183
 - extended, 163
 - ordinary, 163
 - relative, 163
 - extended, 161, 208
 - total, 184, 201
- persistent
 - Betti number, 151, 182
 - homology group, 151, 182
- phenotype, 218
- piecewise linear (see PL), 135
- pixel, 214
- PL
 - critical vertex, 136
 - function, 135, 160
 - Morse
 - function, 137, 202
 - Inequalities, 138
 - regular vertex, 136
- planar graph, 18
- Poincaré
 - duality, 109, 114, 166
 - Duality Theorem, 112, 116
 - homology 3-sphere, 109
 - map, 116
- polygon, 10, 210
- polygonal schema, 29
- polyhedron, 52
- polynomial growth, 184
- positive simplex, 121, 154
- power, 65
 - cell, 65
 - diagram, 65
- priority queue, 43
- projective
 - plane, 29
 - space, 100
- protein, 206
 - docking, 206
 - interaction, 206
- query point, 10
- queue, 198
- quotient topology, 141
- randomized algorithm, 63
- rank of a vector space, 82
- real projective space, 100
- reduced
 - Betti number, 83
 - homology group, 83, 176
 - matrix, 153, 176
- reduction, 176
- Reduction Lemma, 191
- Reeb graph, 141
- regular
 - event, 73
 - point, 127
 - triangulation, 68
 - value, 127
 - vertex, 136
- Reidemeister move, 14
- relative
 - boundary group, 90
 - chain group, 90
 - cycle group, 90
 - homology group, 90, 162, 170
 - persistence diagram, 163
- retract, 58
- retraction, 58
- Riemannian metric, 129
- root architecture, 218
- ru*-decomposition, 177
- saddle, 128
- Schönflies Theorem, 10
- segmentation, 214
 - clock, 199
- separation, 4
- set system, 6
 - convex, 57
- shelling, 24
- short exact sequence, 93
 - of chain complexes, 95
- shortest
 - augmenting path, 193
 - path, 196
- signed distance, 45

- simple
 - closed curve, 9
 - PL critical vertex, 136
- simplex, 51
- simplicial
 - approximation, 56
 - complex, 52
 - homeomorphism, 54
 - map, 54
- Simplicial Approximation Theorem, 56
- simplification, 42, 202, 216
- skeleton, 52
- smallest enclosing ball, 60
- Smith normal form, 87
- smooth, 37
 - function, 126
 - structure, 37
- Snake Lemma, 96
- somite, 199
- space
 - carving, 222
 - curve, 17
- Space Carving Algorithm, 222
- spanning tree, 4
- sparse matrix, 156
- Spectral Sequence
 - Algorithm, 168
 - Theorem, 171
- spectral sequence, 166
- speed, 186
- Sperner Lemma, 101
- sphere, 99
- splay tree, 145
- squared distance, 45
- Stability Theorem
 - for Filtrations, 182
 - for Lipschitz Functions, 185
 - for Tame Functions, 183, 190
 - for Total Persistence, 204
- stable manifold, 131
- standard simplex, 62
- star, 52
 - condition, 55
- Steenrod Five Lemma, 102
- stereographic projection, 64
- Stereographic Projection Lemma, 64
- straight-line
 - embedding, 18
 - homotopy, 181
- strand, 14
 - RNA, 200
- strictly convex combination mapping, 20
- subcomplex, 52
 - full, 52
- subdivision, 54
- sublevel set, 125, 162, 182
- subspace topology, 5
- superlevel set, 162
- surface, 27, 159
- molecular, 208
- surgery, 209
- suspension, 109
- sweeping, 180
- switch, 178
- symbolic perturbation, 13
- symmetry, 163
 - group, 33
- tame, 183
- tangent space, 126
- term in spectral sequence, 169
- Thiessen polygons, 68
- tip, 220
- topological
 - equivalence, 9
 - Morse function, 74
 - space, 5
 - type, 44
- topology, 4
- torus, 32, 102, 116, 125
- total
 - curvature, 186
 - mean curvature, 212
 - persistence, 184, 201
 - variation, 201
- Total Curvature Formula, 188
- transposition, 176
- Transposition Lemma, 179
- transversal, 132
- tree, 4
 - trefoil knot, 13
- triangulable, 52
- triangulation, 30, 52
 - coherent, 68
 - Delaunay, 67
 - irreducible, 48
 - of a polygon, 11
 - regular, 68
- tricoloring, 14
- triple point, 42
- trivial
 - knot, 13
 - link, 15
- Tutte's Theorem, 21, 44
- twisting number, 16
- underlying space, 52
- unfolding, 137
- union of balls, 68
- union-find data structure, 8, 122
- unknot, 13
- unlink, 15
- unstable manifold, 131
- up-tree, 8
- upper
 - link, 215
 - star, 165

van der Waals
 force, 207
 sphere, 208
vector
 field, 128
 space, 93
velocity vector, 186
vertebrate, 199
vertex
 map, 54
 scheme, 53
 set, 52
Vietoris-Rips complex, 61
vine, 181
vineyard, 181
Voronoi
 cell, 65
 weighted, 65, 69
 diagram, 65, 69
 weighted, 65
voxel, 214

Wasserstein distance, 183
Watershed Algorithm, 215
watershed line, 215
weighted
 alpha complex, 69
 Delaunay complex, 68
 graph, 196
 squared distance, 65
 union, 8
 Voronoi
 cell, 65, 69
 diagram, 65
Whitehead link, 25
Whitney umbrella, 38
winding number, 12, 17
writhing number, 16

x-ray crystallography, 207

Combining concepts from topology and algorithms, this book delivers what its title promises: an introduction to the field of computational topology. Starting with motivating problems in both mathematics and computer science and building up from classic topics in geometric and algebraic topology, the third part of the text advances to persistent homology. This point of view is critically important in turning a mostly theoretical field of mathematics into one that is relevant to a multitude of disciplines in the sciences and engineering.

The main approach is the discovery of topology through algorithms. The book is ideal for teaching a graduate or advanced undergraduate course in computational topology, as it develops all the background of both the mathematical and algorithmic aspects of the subject from first principles. Thus the text could serve equally well in a course taught in a mathematics department or computer science department.



ISBN 978-1-4704-6769-2



9 781470 467692

MBK/69.S



For additional information

and updates on this book, visit

www.ams.org/bookpages/mbk-69

