Preface

Modern biology is awash in data. This situation, the result of a technical revolution in high-throughput genomics, promises rapid scientific advances. However, analyzing the data poses unique challenges. Unlike in physics, there is usually no quantitative biological model that can guide investigation and generate precise predictions; often, we do not even know what the relevant quantities are that could capture the essential behavior of the biological system.

In response to the flood of data, the use of clustering algorithms and dimensionality reduction procedures is now ubiquitous. These families of techniques can be regarded as efforts to describe the *shape* of the data set. Although there have been noted successes, such methods provide only crude descriptions of this shape. The power of these tools, as well as their evident limitations, makes it clear that there would be substantial scientific benefit from richer and more robust methods for understanding geometric structure in data.

Algebraic topology is a well-established branch of pure mathematics that studies qualitative descriptors of the shape of geometric objects. Roughly speaking, the goal of algebraic topology is to reduce questions about comparing shapes to questions about comparing algebraic invariants (e.g., numbers), which are typically easier to solve. Moreover, algebraic topology has had a long tradition of employing combinatorial models of geometric objects, *simplicial complexes*, that are well suited to algorithmic computation.

Topological data analysis is a rapidly developing subfield that leverages the tools and outlook of algebraic topology to provide a methodology for analyzing the shape of data sets. The basic strategy is to assign a family of simplicial complexes to a data set; invariants of the complexes integrate information about the shape of the data across different feature scales.

Our aim in this book is to provide a concise introduction to the central ideas and techniques of topological data analysis and to explain in detail a number of specific

Preface

applications to biology. We imagine as our idealized readers a modern quantitative biologist or a graduate student in mathematics with a background in topology or geometry and an interest in applied problems. We have three central goals:

- 1. to equip the modern quantitative biologist with techniques from topological data analysis,
- 2. to direct mathematicians with training in geometry and topology towards problems of interest to biologists, and
- 3. to make it easier for mathematicians and biologists to communicate and collaborate.

These goals pose an expositional challenge, as we expect two quite different audiences with different backgrounds. To address this, we have attempted as much as possible to provide a self-contained introduction to the relevant topics along with abundant and detailed references. We assume that the reader has some familiarity with calculus, linear algebra, elementary probability, and basic statistics.

The first part of this book presents the mathematical background necessary to understand topological data analysis and then provides an overview of techniques in the area. These chapters are intended to be read in order, as each one builds on the previous chapters. The second part of this book consists of a collection of distinct biological applications; each chapter can be read independently.

Acknowledgements

This work grew out of the efforts of many people. We would like to thank Arnold Levine, for his vision, his scientific insights, and his enthusiasm. He created an exceptional creative interdisciplinary environment at the Institute for Advanced Study in Princeton, providing the seeds of many of the ideas discussed in this book. Pablo G. Cámara, Joseph Chan, Kevin Emmett, and Daniel Rosenbloom contributed to several sections in the initial draft of the book. M. Riley Meth made many invaluable corrections and contributions to the second draft of the book. Juan Patino Galindo provided feedback on using genomic data for studying evolutionary processes, and, in particular, helped to write an introduction on different methods to study recombination. We are particularly thankful to Timothy Chu, Oliver Elliott, and M. Riley Meth for using their artistic talents to design the illustrations that enliven the book. William Blumberg and Michael Walfish provided careful readings and helpful comments on previous drafts. Jacqueline Aw, Kyle Bolo, Andrew Chen, Ioan Filip, Chioma Madubata, Patrick Van Nieuwenhuizen, Samuel J. Resnick, Richard T. Wolff, and Sakellarios Zairis proofread different sections of the book. Michael Lesnick and Jun-Hou Fung gave the entire book a very careful reading and made numerous helpful comments correcting errors and

improving the exposition. The authors gratefully acknowledge many interesting discussions with Nils Baas, Gunnar Carlsson, Ben Greenbaum, Gillian Grindstaff, Hossein Khiabanian, Michael Lesnick, Arnold Levine, Michael Mandell, M. Riley Meth, Bud Mishra, Anthea Monod, Sayan Mukherjee, Vladimir Trifonov, Stephen Walker, and Jiguang Wang. In addition, Raúl Rabadán would like to acknowledge many of his collaborators in biology for the time shared, their patience, and the fun solving many problems together: Uttiya Basu, Riccardo Dalla Favera, Adolfo Ferrando, Antonio Iavarone, Anna Lasorella, Tom Maniatis, Do-Hyun Nam, Gustavo Palacios, Teresa Palomero, Laura Pasqualucci, Abbas Rizvi, and Sagi Shapira among many others.

This book was possible in part due to the funding from the Center for Topology and Evolution of Cancer at Columbia University through the National Cancer Institute (U54 CA193313). The Center brings together mathematicians and cancer biologists to solve some interesting problems in cancer. This book was born from many interesting interactions between mathematicians, computational biologists and cancer researchers, where with more or less success, but always with enthusiasm, we have tried to cross the interdisciplinary borders that separate our disciplines. In addition, Raúl Rabadán would like to acknowledge the National Institute of Health grants, R01 CA179044, R01 GM109018, R01 CA185486 and U54 CA209997, and the Convergence program of Stand Up to Cancer together with National Science Foundation. Both authors acknowledge the National Institute of Health grant R01 GM117591. Andrew Blumberg would also like to acknowledge AFOSR research grant FA9550-15-1-0302.