

COMPUTATIONAL PHYLOGENETICS

An Introduction to Designing Methods for Phylogeny Estimation

A comprehensive account of both basic and advanced material in phylogeny estimation, focusing on computational and statistical issues. No background in biology or computer science is assumed, and there is minimal use of mathematical formulas, meaning that students from many disciplines, including biology, computer science, statistics, and applied mathematics, will find the text accessible.

The mathematical and statistical foundations of phylogeny estimation are presented rigorously, following which more advanced material is covered. This includes substantial chapters on multi-locus phylogeny estimation, supertree methods, Markov models of sequence evolution, multiple sequence alignment techniques, and designing methods for large-scale phylogeny estimation. The author provides key analytical techniques to prove theoretical properties about methods, as well as addressing performance in practice for methods for estimating trees. Research problems requiring novel computational methods are also presented, so that graduate students and researchers from varying disciplines will be able to enter the broad and exciting field of computational phylogenetics.

TANDY WARNOW is a Founder Professor of Engineering at the University of Illinois at Urbana-Champaign. Her awards include the National Science Foundation Young Investigator Award (1994), the David and Lucile Packard Foundation Award in Science and Engineering (1996), a Radcliffe Institute for Advanced Study Fellowship (2003), and a John Simon Guggenheim Memorial Foundation Fellowship (2011). She was elected a Fellow of the Association for Computing Machinery (ACM) in 2006, and of the International Society for Computational Biology (ISCB) in 2017.

COMPUTATIONAL PHYLOGENETICS

An Introduction to Designing Methods for Phylogeny Estimation

TANDY WARNOV
University of Illinois, Urbana-Champaign



CAMBRIDGE
UNIVERSITY PRESS

University Printing House, Cambridge CB2 8BS, United Kingdom
One Liberty Plaza, 20th Floor, New York, NY 10006, USA
477 Williamstown Road, Port Melbourne, VIC 3207, Australia
4843/24, 2nd Floor, Ansari Road, Daryaganj, Delhi – 110002, India
79 Anson Road, #06–04/06, Singapore 079906

Cambridge University Press is part of the University of Cambridge.

It furthers the University's mission by disseminating knowledge in the pursuit of education, learning, and research at the highest international levels of excellence.

www.cambridge.org

Information on this title: www.cambridge.org/9781107184718

DOI: 10.1017/9781316882313

© Cambridge University Press 2018

This publication is in copyright. Subject to statutory exception and to the provisions of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published 2018

Printed in the United Kingdom by TJ International Ltd. Padstow Cornwall

A catalogue record for this publication is available from the British Library.

ISBN 978-1-107-18471-8 Hardback

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party internet websites referred to in this publication and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.

Dedicated to Eugene Leighton (Gene) Lawler (1933–1994), my doctoral advisor,
whose enthusiasm and generosity inspired me throughout graduate school,
and who introduced me to the field of computational phylogenetics.

Contents

	<i>Preface</i>	<i>page</i> xiii
	<i>Glossary</i>	xvii
	<i>Notation</i>	xviii
	PART I BASIC TECHNIQUES	1
1	Brief Introduction to Phylogenetic Estimation	3
	1.1 The Cavender–Farris–Neyman Model	4
	1.2 An Analogy: Determining Whether a Coin is Biased Toward Heads or Tails	6
	1.3 Estimating the Cavender–Farris–Neyman Tree	7
	1.4 Some Comments about the CFN Model	16
	1.5 Phylogeny Estimation Methods Used in Practice	16
	1.6 Measuring Error Rates on Simulated Datasets	18
	1.7 Getting Branch Support	20
	1.8 Using Simulations to Understand Methods	20
	1.9 Genome-Scale Evolution	23
	1.10 Designing Methods for Improved Accuracy and Scalability	24
	1.11 Summary	24
	1.12 Review Questions	26
	1.13 Homework Problems	27
2	Trees	29
	2.1 Introduction	29
	2.2 Rooted Trees	29
	2.3 Unrooted Trees	35
	2.4 Constructing the Strict Consensus Tree	41
	2.5 Quantifying Error in Estimated Trees	41
	2.6 The Number of Binary Trees on n Leaves	43
	2.7 Rogue Taxa	43
	2.8 Difficulties in Rooting Trees	44

2.9	Homeomorphic Subtrees	45
2.10	Some Special Trees	45
2.11	Further Reading	46
2.12	Review Questions	47
2.13	Homework Problems	47
3	Constructing Trees from True Subtrees	51
3.1	Introduction	51
3.2	Tree Compatibility	51
3.3	The Algorithm of Aho, Sagiv, Szymanski, and Ullman: Constructing Rooted Trees from Rooted Triples	52
3.4	Constructing Unrooted Binary Trees from Quartet Subtrees	53
3.5	Testing Compatibility of a Set of Trees	56
3.6	Further Reading	57
3.7	Review Questions	58
3.8	Homework Problems	58
4	Constructing Trees from Qualitative Characters	61
4.1	Introduction	61
4.2	Terminology	62
4.3	Tree Construction Based on Maximum Parsimony	63
4.4	Constructing Trees from Compatible Characters	69
4.5	Tree Construction Based on Maximum Compatibility	72
4.6	Treatment of Missing Data	75
4.7	Informative and Uninformative Characters	75
4.8	Further Reading	77
4.9	Review Questions	78
4.10	Homework Problems	78
5	Distance-based Tree Estimation Methods	83
5.1	Introduction	83
5.2	UPGMA	84
5.3	Additive Matrices	86
5.4	Estimating Four-Leaf Trees: The Four Point Method	87
5.5	Quartet-based Methods	89
5.6	Neighbor Joining	91
5.7	Distance-based Methods as Functions	92
5.8	Optimization Problems	94
5.9	Minimum Evolution	95
5.10	The Safety Radius	96
5.11	Comparing Methods	99
5.12	Further Reading	100
5.13	Review Questions	103
5.14	Homework Problems	104

6	Consensus and Agreement Trees	109
	6.1 Introduction	109
	6.2 Consensus Trees	109
	6.3 Agreement Subtrees	116
	6.4 Clustering Sets of Trees	117
	6.5 Further Reading	117
	6.6 Review Questions	118
	6.7 Homework Problems	118
7	Supertrees	121
	7.1 Introduction	121
	7.2 Compatibility Supertrees	123
	7.3 Asymmetric Median Supertrees	123
	7.4 Robinson–Foulds Supertrees	124
	7.5 Matrix Representation with Parsimony	126
	7.6 Matrix Representation with Likelihood	128
	7.7 Quartet-based Supertrees	128
	7.8 The Strict Consensus Merger	132
	7.9 SuperFine: A Meta-Method to Improve Supertree Methods	135
	7.10 Further Reading	139
	7.11 Review Questions	142
	7.12 Homework Problems	142
	PART II MOLECULAR PHYLOGENETICS	143
8	Statistical Gene Tree Estimation Methods	145
	8.1 Introduction to Statistical Estimation in Phylogenetics	145
	8.2 Models of Site Evolution	146
	8.3 Model Selection	151
	8.4 Distance-based Estimation	152
	8.5 Calculating the Probability of a Set of Sequences on a Model Tree	154
	8.6 Maximum Likelihood	157
	8.7 Bayesian Phylogenetics	159
	8.8 Statistical Properties of Maximum Parsimony and Maximum Compatibility	161
	8.9 The Impact of Taxon Sampling on Phylogenetic Estimation	164
	8.10 Estimating Branch Support	165
	8.11 Beyond Statistical Consistency: Sample Complexity	167
	8.12 Absolute Fast Converging Methods	167
	8.13 Heterotachy and the No Common Mechanism Model	170
	8.14 Further Reading	172
	8.15 Review Questions	173

x	<i>Contents</i>	
	8.16 Homework Problems	174
9	Multiple Sequence Alignment	178
	9.1 Introduction	178
	9.2 Evolutionary History and Sequence Alignment	180
	9.3 Computing Differences Between Two Multiple Sequence Alignments	180
	9.4 Edit Distances and How to Compute Them	184
	9.5 Optimization Problems for Multiple Sequence Alignment	190
	9.6 Sequence Profiles	194
	9.7 Profile Hidden Markov Models	198
	9.8 Reference-based Alignments	204
	9.9 Template-based Methods	205
	9.10 Seed Alignment Methods	206
	9.11 Aligning Alignments	207
	9.12 Progressive Alignment	209
	9.13 Consistency	212
	9.14 Weighted Homology Pair Methods	213
	9.15 Divide-and-Conquer Methods	214
	9.16 Co-estimation of Alignments and Trees	215
	9.17 Ensembles of HMMs	220
	9.18 Consensus Alignments	224
	9.19 Discussion	226
	9.20 Further Reading	227
	9.21 Review Questions	231
	9.22 Homework Problems	231
10	Phylogenomics: Constructing Species Phylogenies from Multi-Locus Data	234
	10.1 Introduction	234
	10.2 The Multi-Species Coalescent Model (MSC)	235
	10.3 Using Standard Phylogeny Estimation Methods in the Presence of ILS	238
	10.4 Probabilities of Gene Trees under the MSC	239
	10.5 Coalescent-based Methods for Species Tree Estimation	241
	10.6 Improving Scalability of Coalescent-based Methods	253
	10.7 Species Tree Estimation under Duplication and Loss Models	254
	10.8 Constructing Trees in the Presence of Horizontal Gene Transfer	259
	10.9 Phylogenetic Networks	260
	10.10 Further Reading	268
	10.11 Review Questions	272
	10.12 Homework Problems	272
11	Designing Methods for Large-Scale Phylogeny Estimation	274
	11.1 Introduction	274
	11.2 Standard Approaches	274

<i>Contents</i>		xi
11.3	Introduction to Disk-Covering Methods (DCMs)	279
11.4	DCMs that Use Distance Matrices	282
11.5	Tree-based DCMs	285
11.6	Recursive Decompositions of Triangulated Graphs	288
11.7	Creating Multiple Trees	288
11.8	DACTAL: A General Purpose DCM	289
11.9	Triangulated Graphs	293
11.10	Further Reading	296
11.11	Review Questions	297
11.12	Homework Problems	298
Appendix A	Primer on Biological Data and Evolution	299
Appendix B	Algorithm Design and Analysis	304
Appendix C	Guidelines for Writing Papers About Computational Methods	327
Appendix D	Projects	331
	<i>References</i>	339
	<i>Index</i>	376

Preface

Overview

The evolutionary history of a set of genes, species, or individuals provides a context in which biological questions can be addressed. For this reason, phylogeny estimation is a fundamental step in many biological studies, with many applications throughout biology, such as protein structure and function prediction, analyses of microbiomes, inference of human migrations, etc. In fact, there is a famous saying by Dobzhansky that “Nothing in biology makes sense except in the light of evolution” (Dobzhansky, 1973).

Because phylogenies represent what has happened in the past, they cannot be directly observed but rather must be estimated. Consequently, increasingly sophisticated statistical models of sequence evolution have been developed, and are now used to estimate phylogenetic trees. Indeed, over the last few decades, hundreds of software packages and novel algorithms have been developed for phylogeny estimation, and this influx of computational approaches into phylogenetic estimation has transformed systematics. The availability of sophisticated computational methods, fast computers and high-performance computing (HPC) platforms, and large sequence datasets enabled through DNA sequencing technologies, has led to the expectation that highly accurate large-scale phylogeny estimation, potentially answering open questions about how life evolved on earth, should be achievable.

Yet large-scale phylogeny estimation turns out to be much more difficult than expected, for multiple reasons. First, all the best methods are computationally intensive, and standard techniques do not scale well to large datasets; for example, maximum likelihood phylogeny estimation is NP-hard, so exact solutions cannot be found efficiently (unless $P = NP$), and Bayesian MCMC methods can take a long time to reach stationarity. While massive parallelism can ameliorate these challenges to some extent, it doesn't really address the basic challenge inherent in searching an exponential search space. However, another issue is that the statistical models of sequence evolution that properly address genomic data are substantially more complex than the ones that model individual loci, and methods to estimate genome-scale phylogenies are (relatively speaking) in their infancy compared to methods for single gene phylogenies. Finally, there is a substantial gap between performance as suggested by mathematical theory (which is used to establish guarantees about

methods under statistical models of evolution) and how well the methods actually perform on data – even on data generated under the same statistical models! Indeed, this gap is one of the most interesting things about doing research in computational phylogenetics, because it means that the most impactful research in the area must draw on mathematical theory (especially probability theory and graph theory) as well as on observations from data.

The main goal of this text is to enable researchers (typically graduate students in computer science, applied mathematics, or statistics) to be able to contribute new methods for phylogeny estimation, and in particular to develop methods that are capable of providing improved accuracy for large heterogeneous datasets that are characteristic of the types of inputs that are increasingly of interest in practice. The secondary goal is to enable biologists to understand the methods and their statistical guarantees under these models of evolution, so that they can select appropriate methods for their datasets, and select appropriate datasets given the available methods.

Some of the material in the textbook is fairly mathematical, and presumes undergraduate coursework in discrete mathematics and algorithm design and analysis. However, no background in biology is assumed, and the assumed statistics background is relatively lightweight. While some students without the expected background in computer science may find it difficult to understand some of the proofs, my goal has been to enable all students to understand the theoretical guarantees for phylogeny estimation methods and the statistical models on which they are based, so that they can adequately critique the scientific literature, and also choose methods and datasets that are best able to address the scientific questions they wish to answer.

Outline of the Textbook

Part I provides the “Discrete Mathematics for Phylogenetics” foundations for the textbook; the concepts and mathematics introduced in this part are the building blocks for algorithm design in phylogenetics, especially for developing methods that can scale to large datasets; understanding these concepts makes it possible to understand theoretical guarantees of methods under statistical models of evolution. Chapter 1 introduces the major themes involved in computational phylogenetics, addressing both theory (e.g., statistical consistency under a statistical model of evolution) and performance on both simulated and biological data. This chapter uses the Cavender–Farris–Neyman model of binary sequence evolution since understanding issues in analyzing data generated by this very simple model is helpful to understanding statistical estimation under the commonly used models of molecular sequence evolution. Chapter 2 introduces trees as graph-theoretic objects, and presents different representations of trees that will be useful for method development. Chapters 3, 4, and 5 present different types of methods for phylogenetic tree estimation (based on combining subtrees, using character data, or using distances, respectively). Chapter 6 presents methods for analyzing sets of trees, each on the same set of taxa, and for computing consensus trees and agreement subtrees; it also discusses how these methods are used to estimate support for different phylogenetic hypotheses. Chapter 7 examines the

topic of supertree estimation, where the input is a set of trees on overlapping sets of taxa and the objective is a tree on the full set of taxa. Supertree methods are of interest in their own right and also because they are key algorithmic ingredients in divide-and-conquer methods, a topic we return to in Chapter 11.

Part II of the textbook is concerned with molecular phylogenetics. Chapter 8 presents commonly used statistical models of molecular sequence evolution and statistical methods for phylogeny estimation under these models. However, standard sequence evolution models do not include events such as insertions, deletions, and duplications, which can change the sequence length. These are very common processes, so biological sequences are usually of different lengths and must first be put into a *multiple sequence alignment* before they can be analyzed using phylogeny estimation methods; the subject of how to compute a multiple sequence alignment is covered in Chapter 9. Constructing a species tree or even a phylogenetic network from different gene trees in the presence of gene tree heterogeneity due to incomplete lineage sorting, gene duplication and loss, horizontal gene transfer, or hybridization is a fascinating research area that we present in Chapter 10. We end with Chapter 11, which addresses method development for estimating trees on large datasets. Large-scale phylogeny estimation is increasingly important since nearly all good approaches to phylogeny and multiple sequence alignment estimation are computationally intensive (either heuristics for NP-hard optimization problems or Bayesian methods), and many large datasets are being assembled that cannot be accurately analyzed using existing methods.

Each chapter ends with a set of review questions and homework problems. The review questions are easy to answer and do not require any significant problem solving or calculation. The homework problems are largely pen and paper problems that reinforce the mathematical content of the text.

The textbook comes with four appendices. Appendix A provides an introduction to biological evolution and data; the textbook can be read without it, but the reader who wishes to analyze biological data will benefit from this material. Appendix B provides an introduction to algorithm design and analysis; this material is not necessary for students with undergraduate computer science backgrounds, but may be a helpful introduction for students without this background. Appendix C provides some guidelines about how to write papers that introduce new methods or evaluate existing methods. Appendix D provides computational projects ranging from short term (i.e., a few days) to research projects that could lead to publications. In fact, several of the final projects for my Computational Phylogenetics courses have grown into journal publications (e.g., Bayzid et al. (2014); Zimmermann et al. (2014); Davidson et al. (2015); Chou et al. (2015); Nute and Warnow (2016)).

I wish to thank my editor, David Tranah, for his detailed and insightful comments on the many earlier versions of the text. I also wish to thank my students, colleagues, and family members who gave helpful criticism, including Edward Braun, Sarah Christensen, Steve Evans, Dan Gusfield, Joseph Herman, Ally Kaminsky, Laura Kubatko, Ari Löytynoja, Siavash Mirarab, Erin Molloy, Luay Nakhleh, Mike Nute, David Posada, Bhanu Renukuntla, Ehsan Saleh, Erfan Sayyari, Kimmen Sjölander, Travis Wheeler, and Mark Wilkinson. The several anonymous reviewers also gave very useful comments.

The images of the Monterey Cypress tree on the front and back covers are in honor of the CIPRES project (www.phylo.org), an NSF-funded project for phylogenetic research that I co-lead with Bernard Moret from 2003–2010. Many of the algorithmic advances discussed in the text came out of research supported by CIPRES.

Glossary

afc: Absolute fast-converging methods

ASSU: The algorithm by Aho, Sagiv, Szymanski, and Ullman for determining if a set of rooted triplet trees is compatible, and constructing the compatibility tree if it exists

centroid edge: An edge in a tree T whose deletion defines a decomposition of the leafset into two parts that is as close to balanced as possible

c-gene: A region within a set of genomes that is recombination-free

GTR: General Time Reversible model

GTR-GAMMA: The GTR model with gamma-distributed rates across sites

HMM: Hidden Markov model

homologous: Two sequences are homologous if they have descended from a common ancestor

homoplasy: Evolution with back-mutation or parallel evolution

indels: Insertions and deletions

JC69: Jukes–Cantor model

K2P: Kimura 2-parameter model

MRCA: Most recent common ancestor

MSC: Multi-species coalescent model

NP: The class of decision (i.e., yes/no) problems for which the yes-instances can be verified in polynomial time

NP-hard: A problem that is at least as difficult as the hardest problems in the class NP

NP-complete: A decision problem that is NP-hard and also in the class NP

polytomy: A node in an unrooted tree of degree greater than three, or a node in a rooted tree with more than two children

Notation

- λ : The empty string
- $ab|cd$: Quartet tree on leafset a, b, c, d with one internal edge separating a, b from c, d
- $(a, (b, c))$: Rooted tree on three leaves a, b, c in which b and c are siblings
- $Clades(T)$: The set of clades of a rooted tree T , where a clade is the set of leaves below some internal node in T
- $C(T)$: The set of bipartitions on the leafset induced by edge deletions in a tree T
- $C_I(T)$: The set of non-trivial bipartitions on the leafset induced by deletions of internal edges in a tree T
- $\mathcal{L}(T)$: The set of leaves of a tree T
- $L_\infty(M, M')$: For matrices \mathbf{M} and \mathbf{M}' with the same dimensions, this is $\max_{ij} |M_{ij} - M'_{ij}|$
- $M[i, j]$: For matrix \mathbf{M} , this is the entry in row i and column j . This is also denoted by M_{ij} .
- $MP(T, M)$: The maximum parsimony score of a tree T given the character matrix \mathbf{M}
- $Q(T)$: The set of homeomorphic unrooted quartet trees induced by T on its leafset
- $Q_r(T)$: The set of unrooted fully resolved (i.e., binary) quartet trees in $Q(T)$
- $|S|$: The number of elements in the set S
- $S \setminus S'$: The set $\{x : x \in S \text{ and } x \notin S'\}$ (i.e., the elements of S that are not in S')
- $T|X$: The subtree of T induced on leafset X , with nodes of degree two suppressed
- T_u : The unrooted tree obtained by suppressing the root for T