

Markov Chain Monte Carlo and its Applications to Phylogenetic Tree Construction

by

Marta Magdalena Luczynska

Bachelor of Science, Computer Science and Engineering, M.I.T., 2006
Submitted to the Department of Electrical Engineering and Computer
Science

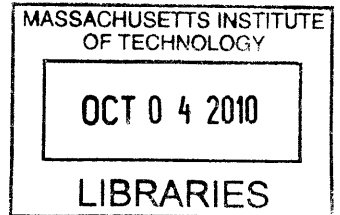
in Partial Fulfillment of the Requirements for the Degree of
Master of Engineering in Electrical Engineering and Computer Science
at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

September 19, 2007

©2007 Massachusetts Institute of Technology
All rights reserved.

ARCHIVES



Author
Department of Electrical Engineering and Computer Science
August 21, 2007

Certified by... ..
Keh-Ping Dunn
Group Leader, Group 32, MIT Lincoln Laboratory
VI-A Company Thesis Supervisor

Certified by.....
Manolis Kellis
Assistant Professor, Distinguished Alumnus (1964) Career
Development
M.I.T. Thesis Supervisor

Accepted by
Arthur C. Smith
Professor of Electrical Engineering
Chairman, Department Committee on Graduate Theses

Markov Chain Monte Carlo and its Applications to Phylogenetic Tree Construction

by

Marta Magdalena Luczynska

Submitted to the
Department of Electrical Engineering and Computer Science
August 21, 2007

In Partial Fulfillment of the Requirements for the Degree of
Master of Engineering in Electrical Engineering and Computer Science

Abstract

This thesis addresses the application of Bayesian methods to problems in phylogenetics. Specifically, we focus on using genetic data to estimate phylogenetic trees representing the evolutionary history of genes and species. Knowledge of this common ancestry has implications for the identification of functions and properties of genes, the effect of mutations and their roles in particular diseases, and other diverse aspects of the biology of cells. Improved algorithms for phylogenetic inference should increase our potential for understanding biological organisms while remaining computationally efficient. To this end, we formulate a novel Bayesian model for phylogenetic tree construction based on recent studies that incorporates known information about the evolutionary history of the species, referred to as the species phylogeny, in a statistically rigorous way. In addition, we develop an inference algorithm for this model based on a Markov chain Monte Carlo method in order to overcome the computational complexity inherent in the problem. Initial results show potential advantages over methods for phylogenetic tree estimation that do not make use of the species phylogeny.

Thesis Supervisor: Keh-Ping Dunn
Title: Group Leader, Group 32, MIT Lincoln Laboratory
VI-A Company Thesis Supervisor

Thesis Supervisor: Manolis Kellis
Title: Assistant Professor, Distinguished Alumnus (1964) Career Development
M.I.T. Thesis Supervisor

Acknowledgments

There have been several people who have provided much support throughout the process of writing this thesis. First, I'd like to thank Keh-Ping Dunn and Manolis Kellis for being such wonderful thesis co-advisors, letting me have a great deal of independence yet always being available to offer help and advice. Next, I'd like to thank Matt Rasmussen for all the time he spent helping me develop the new model and algorithm as well as discussing many of the technical details of the thesis. I'd also like to thank my co-workers at MIT Lincoln Laboratory, including David Choi, Martin Tobias, Pavel Gorelik, Dana Sinno, and Sung-Hyun Son, for taking the time to answer my questions and for offering a great deal of good advice. Lastly, I'd like to thank Dan Rudoy for everything he's done - pushing me forward, being honest, providing moral support, editing chapters time and again, and being someone I could always ask questions (there was an endless number of these) and turn to for help. He seems to have an endless supply of patience. Without him, this thesis would not have been nearly as technically rigorous and none of the equations would have been punctuated.

Contents

1	Introduction	13
2	Markov Chain Monte Carlo	15
2.1	Notation	15
2.2	Introduction	15
2.3	Standard Sampling Methods	18
2.3.1	Inverse Transform Sampling	18
2.3.2	Rejection Sampling	19
2.3.3	Importance Sampling	21
2.4	Discrete Markov Chains	22
2.5	Markov Chain Monte Carlo	26
2.5.1	Metropolis-Hastings Algorithm	26
2.5.2	Gibbs Sampler	29
2.6	Convergence Diagnostics	30
2.7	Practical Considerations	34
2.7.1	One Long Chain vs. Several Shorter Chains	34
2.7.2	Optimizing the Acceptance Rate	35
2.7.3	Burn-in and Choosing a Starting Value	36
2.7.4	Convergence Acceleration	37
2.8	Example: Sampling from a Bivariate Normal Distribution	38
3	Phylogenetics	43
3.1	Genes and their Evolution	43

3.2	Phylogenetic Tree Construction	46
3.2.1	Evolutionary Models	48
3.2.2	Reconciliation	49
3.2.3	Distance-based Approaches	50
3.2.4	Model-based Approaches	54
3.2.5	Tree Proposal Methods	56
4	Phylogenetic Tree Construction Using Markov Chain Monte Carlo	59
4.1	Motivation and Model Description	59
4.1.1	Model Parameters	60
4.1.2	Model Specification	62
4.2	Inference Algorithm	64
4.2.1	Motivation for Developing an MCMC-based Inference Algorithm	64
4.2.2	Algorithm Description	65
4.3	Basic Implementation and Results	67
5	Estimating Gene Phylogenies when Species Phylogeny is Known	71
5.1	New Phylogenetic Tree Model	72
5.1.1	Model Parameters	72
5.1.2	Model Specification	74
5.2	Inference Algorithm	78
5.2.1	Motivation for Developing an MCMC-based Inference Algorithm	79
5.2.2	Algorithm Description	79
5.2.3	Computing the Acceptance Probability	82
5.3	Implementation and Results	84
6	Conclusion	85
A	Proofs	87
A.1	Ergodicity Theorem for Discrete-Time Finite State Space Markov Chains	87
A.2	Counting Distinct Labeled Rooted Bifurcating Trees	89

List of Figures

2-1	Markov chains and their properties	24
2-2	Diagnostic plots for simple MCMC algorithm for sampling from a bi-variate normal distribution	41
3-1	Examples of a rooted and an unrooted bifurcating phylogenetic tree .	47
3-2	Swapping procedure in nearest-neighbor interchange algorithm	57
4-1	Distribution on the set of possible topologies	68
4-2	Results of run of Li <i>et al.</i> 's algorithm	69
5-1	Two different complex reconciliations	75
5-2	Reconciliation of a gene tree	77
A-1	Rooted bifurcating tree with two external nodes	90
A-2	Two ways of adding a new external node to a rooted bifurcating tree	90

List of Tables

2.1	Results of MCMC simulations from a bivariate Gaussian	42
-----	---	----

Chapter 1

Introduction

With the advancements in genome sequencing, we have a vast amount of genetic data at our disposal. The question that arises is how to extract meaningful biological information from this data. One way of extracting such information is through the use of phylogenetic models, which attempt to statistically characterize the evolutionary ties that exist among organisms. One popular class of such models is comprised of phylogenetic trees, which enable us to study organisms' divergence times, and can be used in tasks such as reconstructing ancient proteins, studying the function and properties of genes, aiding in vaccine preparation, identifying mutations associated with a particular disease, and determining the identity of new pathogens [1].

There are several algorithms that exist for phylogenetic tree construction. The most popular examples of distance-based approaches include neighbor-joining, minimum evolution, and parsimony. The other category of approaches is model-based, and includes the maximum likelihood and Bayesian methods. Recently, Bayesian approaches, bolstered by the application of Markov chain Monte Carlo algorithms for inference, have gained popularity as they allow for the incorporation of complex sequence evolution models, produce tree estimates and quantify their uncertainty.

We have considered one such method in detail, the work of Li *et al.* [2]. Their Markov chain Monte Carlo approach to phylogenetic tree construction focuses on reconstructing the evolutionary history of a set of genes. We suspected that incorporating additional information into their model, based on the work of Rasmussen [26],

would improve algorithm performance. This fact motivated us to develop a new model which incorporates the known evolutionary history of the species, known as the species phylogeny. Based on this model, we have developed a new algorithm for phylogenetic tree construction that shows promise for outperforming many of the standard algorithms in existence.

The contributions of this work consist of the elicitation of a new Bayesian model which allows for the incorporation of the species phylogeny into phylogenetic tree construction in a statistically rigorous way. Based on this model, we have developed a new algorithm for phylogenetic tree construction that uses the known species phylogeny to aid the search in the possible space of evolutionary trees. We have also compiled a stand-alone introduction to Markov chain Monte Carlo methods focusing on their biological applications.

The rest of this thesis is structured as follows: Chapter 2 is a stand-alone introduction to Markov chain Monte Carlo methods, which allow one to draw samples from a complex distribution (e.g., the distribution on possible trees in phylogenetic tree construction), and are the basis of the new algorithm we have developed. Chapter 3 provides an overview of the field of phylogenetics, including a review of biological terms necessary to understand the rest of the thesis and a description of phylogenetic tree construction along with the evolutionary models that describe how genes evolve over time. In this chapter, we also summarize the common approaches to phylogenetic tree construction, including both distance and model-based approaches. In Chapter 4, we describe the work of Li *et al.* [2], which was the basis for the new model we have developed, along with our implementation of their algorithm and a review of the results. In Chapter 5 we present the new model and algorithm we have developed. Finally, Chapter 6 provides a summary of our work and possible directions for future research.

Chapter 2

Markov Chain Monte Carlo

2.1 Notation

Throughout this thesis, we employ the following notation:

- Random variables are denoted by capital characters, e.g., X .
- Samples or instantiations of random variables are denoted by subscripts, so that X_i is the i^{th} sample whose distribution is identical to that of X .
- All vectors are denoted by boldface characters. Thus, \mathbf{X} defines a multi-dimensional random variable. A sample i distributed according to \mathbf{X} is denoted by \mathbf{X}_i .

2.2 Introduction

This chapter aims to serve as a stand-alone introduction to Markov chain Monte Carlo (MCMC) methods, which provide a means for sampling from complex distributions (e.g., ones that are otherwise hard to sample from or have no closed-form). Several standard techniques exist for sampling from distributions, including the inverse transform method [3] as well as rejection and importance sampling [3,4]. We describe these techniques in further detail in the next section; however, we note that they are not applicable in all situations, whereas MCMC may offer a viable alternative.

The need to sample from complex distributions often arises in Bayesian analysis, which consists of the development and application of statistical models to explain patterns in observed data. Let the observed data Y and the model parameters θ be random variables. In a Bayesian setting, the relationship between θ and Y is encapsulated in the joint distribution of θ and Y which we denote by $P(\theta, Y)$. We are interested in computing the posterior probability, $P(\theta|Y)$, which is the probability of the model parameters given the observed data. Using Bayes' rule, we know that the posterior probability is proportional to the product of the likelihood $P(Y|\theta)$, which is the probability of seeing the observed data given the model parameters, and the prior probability of the model parameters $P(\theta)$ as shown by the following formula:

$$P(\theta|Y) = \frac{P(Y|\theta)P(\theta)}{\int P(Y|\theta)P(\theta) d\theta} \propto P(Y|\theta)P(\theta). \quad (2.1)$$

Statistical inference is concerned with obtaining point estimates and other statistical quantities which summarize the information in the posterior distribution. Two point estimates prevalent in engineering applications are the maximum a posteriori (MAP) and the minimum mean-squared error (MMSE) estimates. As the name suggests, the MAP estimate $\hat{\theta}$ is the one for which the estimates' posterior probability is maximized. Computing this quantity analytically requires that the numerator in (2.1) has a closed form. Otherwise, numerical methods (e.g., sampling) are necessary. The MMSE point estimate $\tilde{\theta}$, is an estimate of θ with the smallest mean-squared error. It is well known that $\tilde{\theta}$ can be found by computing the posterior mean (i.e., $\tilde{\theta} = E[P(\theta|Y)]$). The calculation of the MMSE estimate requires the ability to evaluate the integral in the denominator of (2.1), which is a normalizing constant that ensures that the probabilities of all the possible value assignments of the parameters sum to one. The computation of this normalizing constant is often intractable for complex models, motivating the application of sampling methods.

Sampling from the posterior distribution enables us to approximate the point estimates described above when exact expressions are unavailable. Consider, for instance, the MMSE estimate. If we were to draw N random samples $\{\mathbf{X}_i; i =$

$1, 2, \dots, N\}$ distributed according to the posterior distribution of \mathbf{X} , which we denote by $\pi(\mathbf{x})$, we can then approximate the expectation of $\pi(\mathbf{x})$, which we denote by $\boldsymbol{\mu}$, by taking the mean of these samples:

$$E_{\pi}[\mathbf{X}_i] \approx \frac{1}{N} \sum_{i=1}^N \mathbf{X}_i.$$

The Strong Law of Large Numbers guarantees that the sample average converges almost surely to the true mean of $\pi(\mathbf{x})$, that is:

$$P\left(\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbf{X}_i = \boldsymbol{\mu}\right) = 1. \quad (2.2)$$

It is important to point out that the MAP and MMSE point estimates do not provide us with any sort of uncertainty quantification (i.e., measure of confidence) which captures the shape of the posterior distribution. However, when employing sampling methods, we can determine not only the expected value and the variance of our estimator, but also compute confidence intervals and quantiles about these estimates.

Some of the available sampling methods face well-known limitations. The inverse transform and rejection sampling methods are not applicable when the normalizing constant for the posterior distribution is unknown or its computation is intractable, as is often the case. Importance sampling does not require the knowledge of normalizing constants, but how well it works depends on the quality of the proposal distribution, which is an estimate of the posterior distribution. Unfortunately, good proposal distributions are not always easy to find (e.g., little information may be available about the modes of the posterior distribution). On the other hand, MCMC methods can be applied when the normalizing constant is unknown or hard to compute and depend on a proposal distribution which does not need to be an estimate of the posterior distribution. Due to their applicability, MCMC methods have enjoyed success in a wide range of fields including statistics, physics, bioinformatics, and genetics [3]. These methods have been used in Bayesian image analysis [5], modeling HIV progression [6], linkage analysis in gene families [7], reconstructing phylogenetic

trees [2], and regression analysis [8] to name but a few applications.

The rest of the chapter is structured as follows: Section 2.3 provides an overview of three standard techniques for drawing samples from a distribution of interest, which we mentioned above, namely, the inverse transform method, rejection sampling and importance sampling. The Markov chain theory background for MCMC is presented in Section 2.4. Section 2.5 describes two popular MCMC methods — the Metropolis-Hastings and the Gibbs samplers. Techniques for assessing and accelerating convergence of MCMC algorithms are outlined in Section 2.6. Practical considerations, such as selection of the starting value and deciding how many MCMC chains should be run, are discussed in Section 2.7. Finally, in Section 2.8 we provide a concrete example that ties together the concepts covered throughout the chapter.

2.3 Standard Sampling Methods

This section provides an overview of three standard techniques for drawing samples from a distribution – the inverse transform method, rejection sampling, and importance sampling. We highlight the cases when these techniques are not applicable but the Markov chain Monte Carlo approach holds promise.

2.3.1 Inverse Transform Sampling

The simplest way of generating random samples from a distribution is through inverse transform sampling. If we let \mathbf{X} be a d -dimensional random variable (i.e., $\mathbf{X} \in \mathbb{R}^d$) with cumulative distribution function (cdf) $F(\mathbf{x})$, then we can generate N samples distributed according to $F(\mathbf{x})$ by sampling the uniform distribution and evaluating the inverse cdf $F^{-1}(\mathbf{x})$ at these samples. The pseudocode for this algorithm is shown in Algorithm 1. To see that Algorithm 1 produces N samples distributed according to $F(\mathbf{x})$, we prove that $F^{-1}(\mathbf{U})$ (a function that evaluates the inverse cdf at values sampled according to the uniform distribution) has distribution F . To see this, note

that F^{-1} is right-continuous, and hence satisfies the following two inequalities:

$$F(F^{-1}(\mathbf{u})) \geq \mathbf{u}, \quad \forall \mathbf{u} \in [0, 1]^d,$$

$$F^{-1}(F(\mathbf{x})) \leq \mathbf{x}, \quad \forall \mathbf{x} \in F^{-1}([0, 1]^d).$$

Therefore,

$$P(F^{-1}(\mathbf{U}) \leq \mathbf{x}) = P(\mathbf{U} \leq F(\mathbf{x})) = F(\mathbf{x}),$$

which proves that the random variable $F^{-1}(\mathbf{U})$ is distributed according to $F(\mathbf{x})$.

Algorithm 1: Inverse Transform Sampling

Input: None

Result: N samples $\{\mathbf{X}_i; i = 1, 2, \dots, N\}$ distributed according to $F(\mathbf{x})$

begin

for $i = 1$ to N **do**

 Draw $\mathbf{U} \sim \mathcal{U}[0, 1]^d$

$\mathbf{X}_i = F^{-1}(\mathbf{U})$

end

This method is not applicable if the inverse cdf cannot be computed. Even when the inverse cdf exists, evaluating $F^{-1}(\mathbf{U})$ can be computationally expensive depending on the particular distribution of interest and the dimensionality of the problem.

2.3.2 Rejection Sampling

In rejection sampling, drawing samples from a distribution denoted by $\pi(\mathbf{x})$ requires finding an envelope function $h(\mathbf{x})$, which is easy to sample from, such that:

$$Mh(\mathbf{x}) \geq \pi(\mathbf{x}), \quad \forall \mathbf{x},$$

where $M \geq 0$. Samples are drawn from $h(\mathbf{x})$, and each sampled point \mathbf{X}_i is accepted with probability $\pi(\mathbf{X}_i)/Mh(\mathbf{X}_i)$. The rejected points are discarded. Sampling continues until a desired number of points has been accepted. The pseudocode for the algorithm is shown in Algorithm 2.

Algorithm 2: Rejection Sampling

Input: $h(\mathbf{x})$ and $\pi(\mathbf{x})$

Result: N samples $\{\mathbf{X}_i; i = 1, 2, \dots, N\}$ distributed according to $\pi(\mathbf{x})$

begin

for $i = 1$ *to* N **do**

repeat

 Draw a sample point \mathbf{Y} according to $h(\mathbf{x})$

 Draw a sample \mathbf{U} according to the uniform distribution on $[0, 1]^d$

if $\mathbf{U} \leq \pi(\mathbf{Y})/Mh(\mathbf{Y})$ **then**

 Accept \mathbf{Y}

until \mathbf{Y} *is accepted*

 Set $\mathbf{X}_i = \mathbf{Y}$

end

We will show why the accepted samples are distributed according to the density proportional to $\pi(\mathbf{x})$ for the 1-dimensional case. The generalization to higher dimensions is straightforward. Consider the following argument:

$$P(X \leq y \text{ and } X \text{ is accepted}) = \int_{-\infty}^y \frac{\pi(x)h(x)}{h(x)} dx = \int_{-\infty}^y \pi(x) dx.$$

Therefore,

$$P(X \text{ is accepted}) = \int_{-\infty}^{\infty} \pi(x) dx = 1,$$

and

$$P(X \leq y | X \text{ is accepted}) = \frac{\int_{-\infty}^y \pi(x) dx}{1} = P(X \leq y),$$

which proves that the accepted values are distributed according to $\pi(x)$.

The time taken to generate N accepted samples depends on the quality of the envelope function $h(\mathbf{x})$. If $Mh(\mathbf{x})$ is much larger than $\pi(\mathbf{x})$, then many sampled points \mathbf{Y} will be rejected for every one that is kept because the ratio $\frac{\pi(\mathbf{Y})}{Mh(\mathbf{Y})}$ will be small. Thus, $Mh(\mathbf{x})$ should be chosen so that it is both close to $\pi(\mathbf{x})$ (to avoid a high number of rejected samples) and is easy to evaluate and simulate from (to avoid a high computational cost since many samples are rejected). Another limitation of rejection sampling is that, like the inverse transform method, it can only be applied when the normalizing constant for $\pi(\mathbf{x})$ is known.

2.3.3 Importance Sampling

Importance sampling is another method for sampling from $\pi(\mathbf{x})$ with the aid of an easy-to-sample-from distribution $q(\mathbf{x})$. We draw each sample according to $q(\mathbf{x})$ and then weigh it by $\frac{\pi(\mathbf{x})}{q(\mathbf{x})}$. The pseudocode is shown in Algorithm 3. Once we have the N samples and weights, we can approximate the expectation of any function $f(\mathbf{x})$ of interest with respect to $\pi(\mathbf{x})$ by taking a weighted average of the samples as shown by the equation below:

$$E_{\pi}[f(\mathbf{x})] \approx \frac{1}{N} \sum_{i=1}^N \frac{\pi(\mathbf{X}_i) f(\mathbf{X}_i)}{q(\mathbf{X}_i)} = \frac{1}{N} \sum_{i=1}^N w_i f(\mathbf{X}_i).$$

The Strong Law of Large Numbers guarantees that this approximation converges to the desired quantity as the number of samples increases (i.e., as $N \rightarrow \infty$). The

Algorithm 3: Importance Sampling

Input: $q(\mathbf{x})$ and $\pi(\mathbf{x})$

Result: N samples $\{\mathbf{X}_i; i = 1, 2, \dots, N\}$ distributed according to $\pi(\mathbf{x})$

begin

 Draw $M \geq N$ samples from $q(\mathbf{x})$

 Assign to the i^{th} sample, \mathbf{X}_i , the weight $w_i = \frac{\pi(\mathbf{X}_i)}{q(\mathbf{X}_i)}$

end

quality of the samples generated by importance sampling largely depends on the particular form of the proposal distribution $q(\mathbf{x})$. The proposal distribution should have the same or greater support than that of $\pi(\mathbf{x})$ in order to guarantee that the samples generated can accurately represent the distribution specified by $\pi(\mathbf{x})$. However, even with the same support, samples generated from $q(\mathbf{x})$ will be ones from high probability regions of $q(\mathbf{x})$. If these regions do not overlap with the high probability regions of $\pi(\mathbf{x})$, many samples must be generated before a set of samples is produced which yields a good estimate of the desired expectation, rendering this approach inefficient.

2.4 Discrete Markov Chains

This section provides the Markov chain theory required to understand the Markov chain Monte Carlo algorithms that form the basis for this thesis. Briefly, a Markov chain is a stochastic process in which the effect of the past on the future is summarized by a **state** that changes over time according to a pre-specified probability distribution. Here we consider only discrete-time finite state space Markov chains. These are chains in which the process is defined at integer values of time, denoted by n , and the state of the chain belongs to a finite set S of possible states, called the **state space**. For background on continuous-time and general state space Markov chains, we refer the reader to [3].

We denote the state of the Markov chain at time n by X_n , and assume that X_n can be one of the m possible states, where $m \in \mathbb{Z}_+$, i.e., $S = \{1, \dots, m\}$. At each time instant n , the **transition probabilities** p_{ij} specify the probability of moving to a new state j given that the Markov chain is currently in state i . In particular, we have that:

$$p_{ij} = P(X_{n+1} = j | X_n = i) \quad \forall i, j \in S. \quad (2.3)$$

A Markov chain must also satisfy the **Markov property**, which requires that the value of the next state only depends on the value of the current state. This means that for all times n , all states $i, j \in S$, and all possible previous state values $\{i_0, \dots, i_{n-1}\}$ we have that:

$$P(X_{n+1} = j | X_n = i, X_{n-1} = i_{n-1}, \dots, X_0 = i_0) = P(X_{n+1} = j | X_n = i) = p_{ij}. \quad (2.4)$$

We can summarize the transition probabilities by a **transition probability matrix** \mathbf{P} , which at row i and column j contains the value p_{ij} . These properties are formalized in the following definition:

Definition 2.4.1. *A discrete-time, finite state space Markov chain is a sequence of random variables X_0, X_1, \dots that*

- *Take on values in a finite state space S*

- Evolve over time according to a transition probability matrix \mathbf{P}
- Satisfy the Markov property defined in Equation (2.4)

The probability of any particular sequence of future states can be computed as:

$$P(X_0 = i_0, X_1 = i_1, \dots, X_n = i_n) = P(X_0 = i_0)p_{i_0i_1}p_{i_1i_2} \cdots p_{i_{n-1}i_n}. \quad (2.5)$$

We may also compute the n -step transition probability $r_{ij}(n)$, which is the probability that, in n time steps, the Markov chain will be in state j , after starting in state i . We can derive this probability as follows:

$$\begin{aligned} r_{ij}(n) &= P(X_n = j | X_0 = i) = \sum_{k=1}^m P(X_n = j | X_{n-1} = k, X_0 = i) P(X_{n-1} = k | X_0 = i) \\ &= \sum_{k=1}^m r_{ik}(n-1)p_{kj}, \quad \text{with } r_{ij}(1) = p_{ij}, \end{aligned}$$

where we used the Markov property to rewrite the summation in the last step.

There are several important properties which describe individual states and the entire Markov chain. Identifying these properties allows us to characterize the overall behavior of the chain (e.g., determine how the chain behaves after many iterations). These properties are summarized in Definitions 2.4.2–2.4.7 below and examples are shown in Figure 2-1.

Definition 2.4.2. A state j is **accessible** from state i if for some n , $r_{ij}(n) > 0$. Intuitively, this means that if we are currently in state i , then the probability we will be in state j at some time in the future is greater than zero.

Definition 2.4.3. A state j **communicates** with state i if state j is accessible from state i and state i is accessible from state j .

Definition 2.4.4. A Markov chain is **irreducible** if all states in the chain communicate with each other. Intuitively, this means that it is possible to get to any state in the chain from any starting state. Otherwise, the Markov chain is **reducible**.

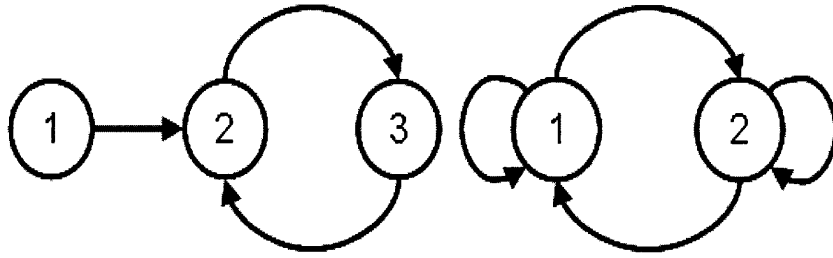


Figure 2-1: Markov chains and their properties. In the Markov chain on the left, state 2 is **accessible** from states 1 and 3, but state 1 is not accessible from state 2 or state 3. State 2 and 3 **communicate**. State 1 is **transient** while states 2 and 3 are **recurrent**. States 2 and 3 are both **periodic** with period 2 while state 1 is **aperiodic**. The entire Markov chain is **reducible** and not ergodic. In the Markov chain on the right, state 1 and state 2 communicate with each other, and are both recurrent and aperiodic. The overall chain is **irreducible** and **ergodic**.

Definition 2.4.5. Let T_j be defined as the return time to state j :

$T_j = \min \{n > 0 : X_n = j | X_0 = j\}$. Then state j is **recurrent** if $P(T_j < \infty) = 1$. A state is called **transient** if it is not recurrent.

Definition 2.4.6. A state j has **period** k , where $k = \gcd\{n : P(X_n = j | X_0 = j)\}$. If $k = 1$, the state is **aperiodic**; otherwise, the state is **periodic** with period k .

Definition 2.4.7. A state is **ergodic** if it is aperiodic and recurrent. Finally, a Markov chain is **ergodic** if all the states in the chain are **ergodic**.

We are often interested in the long-term behavior (i.e., when n is large) of a Markov chain. Denote the probability of being in state j at time n by p_{ij}^n . Then the **limiting probability** π_j , is the probability of being in state j in the limit as the number of time steps goes to infinity, as shown below:

$$\pi_j = \lim_{n \rightarrow \infty} P(X_n = j) = \lim_{n \rightarrow \infty} p_{ij}^n. \quad (2.6)$$

For a discrete-time finite state space Markov chain, aperiodicity and irreducibility are necessary and sufficient conditions for the limiting probabilities to exist and be finite. If the chain is periodic, then some states in the chain can only be visited at certain integer multiples of time so there is no limiting probability. If the chain is reducible, then as time goes on, the chain will eventually enter a state i which does not

communicate with one of the other states j . Since state i and j do not communicate, after this point in time, the chain will never again enter state j . Therefore, the limiting probability of being in state j is 0.

The key question that arises is how to compute each π_j or, written more compactly, $\tilde{\boldsymbol{\pi}} \in \mathbb{R}^{[1 \times m]}$, where $\tilde{\boldsymbol{\pi}} = (\pi_1, \pi_2, \dots, \pi_m)$. First we require another definition:

Definition 2.4.8. *A probability distribution $\boldsymbol{\pi}$ is said to be the **stationary distribution** of the Markov chain with transition matrix \mathbf{P} if it satisfies the following system of equations:*

$$\pi_j = \sum_{k=1}^m \pi_k p_{kj}, \quad \mathbf{1} = \sum_{k=1}^m \pi_k, \quad j = 1, \dots, m \quad (2.7)$$

Equivalently, in matrix notation, $\boldsymbol{\pi} = \boldsymbol{\pi}\mathbf{P}$. We are interested in the conditions that guarantee the equality of $\boldsymbol{\pi}$ and $\tilde{\boldsymbol{\pi}}$. Theorem 2.4.1 states that ergodicity of the Markov chain is a sufficient condition for this to occur.

Theorem 2.4.1. *An ergodic discrete-time finite state space Markov chain converges to a unique stationary distribution $\boldsymbol{\pi}$ (for proof, see Appendix A).*

Theorem 2.4.1 gives us a procedure for determining the limiting distribution of a discrete-time finite state space Markov chain. If we prove that the Markov chain is irreducible and aperiodic (i.e., ergodic), then the vector $\boldsymbol{\pi}$ which satisfies the stationary equations corresponds to the stationary distribution of the Markov chain.

In the next section, we will show that equivalent conditions to those described in Theorem 2.4.1 are satisfied by the Markov chain underlying the MCMC algorithms we discuss. The only thing which we will have to prove is whether or not that Markov chain's stationary distribution is the one we seek. This may be checked by first proving that our chain is ergodic, and then that it satisfies the following **detailed balance equations**:

$$p_{ij}\pi_i = p_{ji}\pi_j, \quad \forall \text{ states } i, j, \quad (2.8)$$

which ensure that the Markov chain's stationary distribution will indeed be $\boldsymbol{\pi}$. To see why this is the case, we can compute the summation over all of the states i on

both sides of the detailed balance equations (2.8) as follows:

$$\sum_{i=1}^m p_{ij} \pi_i = \sum_{i=1}^m p_{ji} \pi_j = \pi_i. \quad (2.9)$$

The middle summation of (2.9) is the marginal distribution of state i while the RHS is the limiting distribution of state i . Therefore we have shown that if an ergodic Markov chain satisfies the detailed balance equations, it has π as its stationary distribution and will eventually converge to it.

2.5 Markov Chain Monte Carlo

The central idea behind the Markov chain Monte Carlo method is to construct a Markov chain with stationary distribution π – the very distribution from which we wish to simulate. We can then draw samples from this distribution of interest by running the Markov chain. The question of how to create such chains is answered constructively. In fact, the two most popular MCMC algorithms, the Metropolis-Hastings and Gibbs samplers, are nothing but explicit recipes for how to construct the Markov chain with the desired properties.

2.5.1 Metropolis-Hastings Algorithm

The Metropolis-Hastings algorithm was developed by Hastings [9] as a generalization of the original method proposed by Metropolis *et al.* [10]. The algorithm begins with some initial guess \mathbf{X}_0 . At iteration i of the algorithm, we sample a candidate point \mathbf{Y}_i from a proposal distribution $q(\mathbf{y}|\mathbf{X}_i)$, which describes how we generate potential new samples given the current sample \mathbf{X}_i . The candidate point is accepted with probability $\alpha(\mathbf{X}_i, \mathbf{Y}_i)$ defined by:

$$\alpha(\mathbf{X}_i, \mathbf{Y}_i) = \min \left(1, \frac{\pi(\mathbf{Y}_i)q(\mathbf{X}_i|\mathbf{Y}_i)}{\pi(\mathbf{X}_i)q(\mathbf{Y}_i|\mathbf{X}_i)} \right). \quad (2.10)$$

The original algorithm developed by Metropolis [10] assumed a symmetric pro-

positional distribution, i.e., $q(\mathbf{X}_i|\mathbf{Y}_i) = q(\mathbf{Y}_i|\mathbf{X}_i)$, in which case the acceptance probability (2.10) simplifies to $\alpha(\mathbf{X}_i, \mathbf{Y}_i) = \min\left(1, \frac{\pi(\mathbf{Y}_i)}{\pi(\mathbf{X}_i)}\right)$. If the candidate point is accepted, it becomes the next state, that is, $\mathbf{X}_{i+1} = \mathbf{Y}_i$. The overall procedure is summarized in Algorithm 4. Intuitively, the Metropolis-Hastings algorithm performs a probabilistic exploration of the state space of \mathbf{X} according to $\pi(\mathbf{x})$. If the new sample proposed has higher probability than \mathbf{X}_i according to $\pi(\mathbf{x})$, then it is always accepted. Otherwise, the sample is kept in proportion to its probability relative to that of the current sample. Hence, as we run the Metropolis-Hastings algorithm, we produce more samples in the higher and fewer samples in the lower probability regions of $\pi(\mathbf{x})$.

We emphasize that the normalizing constant of $\pi(\mathbf{x})$ need not be known since it cancels out in Equation (2.10). If a proposed point is rejected, then the next state remains at the current state, i.e., $\mathbf{X}_{i+1} = \mathbf{X}_i$. Clearly, the samples that we generate are not independent since each new sample is correlated with the sample from the previous iteration through the proposal distribution $q(\mathbf{y}|\mathbf{X}_i)$. In particular, the chain may remain in the same place and hence the new sample will be the same as the current one. Once we have generated our samples, we can estimate the desired expectation of any function of interest with respect to our target density $\pi(\mathbf{x})$ by computing the sample mean.

Algorithm 4: Metropolis-Hastings

Input: $q(\mathbf{y}|\mathbf{x})$ and $\pi(\mathbf{x})$
Result: N samples $\{\mathbf{X}_i; i = 1, 2, \dots, N\}$ that are distributed according to $\pi(\mathbf{x})$ as $N \rightarrow \infty$

```

begin
  Initialize  $\mathbf{X}_0$  to any value in the state space of  $\mathbf{X}$ 
  for  $i = 1$  to  $N$  do
    Sample a point  $\mathbf{Y}_i$  according to  $q(\mathbf{y}|\mathbf{X}_{i-1})$ 
    Sample a uniform random variable  $\mathbf{U} \sim \mathcal{U}[0, 1]$ 
    if  $\mathbf{U} \leq \alpha(\mathbf{X}_i, \mathbf{Y}_i)$  then
      |  $\mathbf{X}_{i+1} = \mathbf{Y}_i$ 
    else
      |  $\mathbf{X}_{i+1} = \mathbf{X}_i$ 
  end
end

```

There are a number of approaches to choosing a proposal distribution $q(\mathbf{y}|\mathbf{X}_i)$.

The first approach, termed random walk by practitioners, involves setting the value of the new sample \mathbf{Y}_i to a perturbed version of the previous sample value, i.e., $\mathbf{Y}_i = \mathbf{X}_i + \mathbf{Z}$ where \mathbf{Z} is a multivariate normal random variable with a tunable covariance matrix Σ . This implies that $\mathbf{Y}_i \sim \mathcal{N}(\mathbf{X}_i, \Sigma)$. The important point to consider when choosing the proposal distribution is the distribution's scale (i.e., Σ^{-1}). A small covariance matrix will result in smaller steps being proposed and hence less of the state space will be explored. A large covariance matrix will result in larger steps being proposed and hence greater exploration of the state space. The second approach is the independence chain sampling approach, in which \mathbf{Y}_i is chosen independently of the current value of \mathbf{X} , i.e., $q(\mathbf{y}|\mathbf{X}_i) = q(\mathbf{y})$. In this case, the acceptance probability in (2.10) can be rewritten as follows:

$$\alpha(\mathbf{X}_i, \mathbf{Y}_i) = \min \left(1, \frac{\pi(\mathbf{Y}_i)q(\mathbf{X}_i)}{\pi(\mathbf{X}_i)q(\mathbf{Y}_i)} \right). \quad (2.11)$$

The key to having the independence chain sampler work well is to choose $q(\mathbf{y})$ to be a good approximation to $\pi(\mathbf{x})$ though slightly heavier-tailed. This is a similar condition to that required for importance sampling to work well. The method is also similar to rejection sampling, except that it produces a sample at every iteration as there are no rejections. If $q(\mathbf{y})$ is chosen to be a good, though slightly heavier-tailed approximation to $\pi(\mathbf{x})$, we will have a lower overall acceptance rate at the benefit of avoiding spending long time intervals in the tails of the target distribution. The samples, as in the random walk approach, are still correlated through the proposal distribution (i.e., they are not independent samples from $\pi(\mathbf{x})$ because we can stay in the same place).

We can also combine several different proposal densities using a mixture approach. As motivation, we note that the independence Metropolis-Hastings algorithm, for example, is prone to missing the finer details of the target distribution $\pi(\mathbf{x})$, especially when the target distribution in question has high dimensionality. This issue can be addressed by utilizing global proposal densities to determine the high probability regions of the state space and local proposals that allow us to explore the finer details

of each of these regions. Our proposal density $q(\mathbf{x})$ in this case will be a weighted sum of N distributions, each which have an invariant distribution of $\pi(\mathbf{x})$. The proposal density will take the following form:

$$q(\mathbf{x}) = \sum_{i=1}^N \alpha_i q_i(\mathbf{x}), \text{ where } \sum_{i=1}^N \alpha_i = 1.$$

For example, we can compose our mixture of an independence proposal density for determining the modes of the target distribution and a random walk proposal density for exploring the state space around each mode.

2.5.2 Gibbs Sampler

Algorithm 5: Gibbs Sampler

Input: $q(\mathbf{y}|\mathbf{x})$ and $\pi(\mathbf{x})$

Result: N samples $\{\mathbf{X}_i; i = 1, 2, \dots, N\}$ that are distributed according to $\pi(\mathbf{x})$ as $N \rightarrow \infty$

begin

 Initialize \mathbf{X}_0 to any value in the state space of \mathbf{X}

for $i = 1$ to N **do**

for each dimension j of \mathbf{X}_i **do**

 Sample \mathbf{X}_{i+1}^j according to $q_j(\mathbf{X}_{i+1}|\mathbf{X}_i)$

end

Another very popular Markov chain Monte Carlo algorithm is the Gibbs sampler [11]. Though it is typically described separately from the Metropolis-Hastings algorithm, it really is a special case of this scheme. The Gibbs sampler is applicable in settings where the dimension d of the density of interest is greater than one and the distributions of each dimension conditioned on the rest are available. The main distinction from the generic Metropolis scheme is in the form of the proposal distribution, which updates each component j of the i^{th} sample of \mathbf{X} , denoted by \mathbf{X}_i^j , in turn. The proposal distribution for updating the j^{th} dimension of the new sample \mathbf{X}_{i+1} is given by the following conditional distribution:

$$q_j(\mathbf{X}_{i+1}|\mathbf{X}_i) = \pi_j(\mathbf{X}_{i+1}^j|\mathbf{X}_{i+1}^1, \mathbf{X}_{i+1}^2, \dots, \mathbf{X}_{i+1}^{j-1}, \mathbf{X}_i^{j+1}, \mathbf{X}_i^{j+2}, \dots, \mathbf{X}_i^d).$$

Therefore, the proposal distribution can be written as follows:

$$q(\mathbf{X}_{i+1}|\mathbf{X}_i) = \prod_{j=1}^d q_j(\mathbf{X}_{i+1}^j|\mathbf{X}_i). \quad (2.12)$$

The acceptance probability for this proposal distribution is 1. We can see this easily since each of the terms in the product in Equation (2.12) takes the following form:

$$\begin{aligned} \frac{\pi(\mathbf{X}_{i+1})q_j(\mathbf{X}_i|\mathbf{X}_{i+1})}{\pi(\mathbf{X}_i)q_j(\mathbf{X}_{i+1}|\mathbf{X}_i)} &= \frac{\pi(\mathbf{X}_{i+1})\pi_j(\mathbf{X}_i^j|\mathbf{X}_i^1, \mathbf{X}_i^2, \dots, \mathbf{X}_i^{j-1}, \mathbf{X}_i^{j+1}, \mathbf{X}_i^{j+2}, \dots, \mathbf{X}_i^d)}{\pi(\mathbf{X}_i)\pi_j(\mathbf{X}_{i+1}^j|\mathbf{X}_i^1, \mathbf{X}_i^2, \dots, \mathbf{X}_i^{j-1}, \mathbf{X}_i^{j+1}, \mathbf{X}_i^{j+2}, \dots, \mathbf{X}_i^d)} \\ &= \frac{\pi_j(\mathbf{X}_{i+1}^j|\mathbf{X}_i^1, \mathbf{X}_i^2, \dots, \mathbf{X}_i^{j-1}, \mathbf{X}_i^{j+1}, \mathbf{X}_i^{j+2}, \dots, \mathbf{X}_i^d)}{\pi_j(\mathbf{X}_i^j|\mathbf{X}_i^1, \mathbf{X}_i^2, \dots, \mathbf{X}_i^{j-1}, \mathbf{X}_i^{j+1}, \mathbf{X}_i^{j+2}, \dots, \mathbf{X}_i^d)} * \\ &= \frac{\pi_j(\mathbf{X}_i^j|\mathbf{X}_i^1, \mathbf{X}_i^2, \dots, \mathbf{X}_i^{j-1}, \mathbf{X}_i^{j+1}, \mathbf{X}_i^{j+2}, \dots, \mathbf{X}_i^d)}{\pi_j(\mathbf{X}_{i+1}^j|\mathbf{X}_i^1, \mathbf{X}_i^2, \dots, \mathbf{X}_i^{j-1}, \mathbf{X}_i^{j+1}, \mathbf{X}_i^{j+2}, \dots, \mathbf{X}_i^d)} = 1. \end{aligned}$$

Therefore,

$$\alpha(\mathbf{X}_{i+1}, \mathbf{X}_i) = \min \left(1, \prod_{j=1}^d \frac{\pi(\mathbf{X}_{i+1})q_j(\mathbf{X}_i|\mathbf{X}_{i+1})}{\pi(\mathbf{X}_i)q_j(\mathbf{X}_{i+1}|\mathbf{X}_i)} \right) = \min(1, 1^d) = 1.$$

Note that since the acceptance probability is always 1, no samples are discarded when running the Gibbs sampler. The pseudocode is shown in Algorithm 5. Under reasonably general conditions (see [12]), as the number of samples goes to infinity, the distribution of the samples generated by this method converges to the target distribution $\pi(\mathbf{x})$. For a more detailed tutorial on the Gibbs sampler, including several examples, we refer the reader to [12]. The Gibbs sampler will be the underlying method in the algorithms presented in Chapter 4 and Chapter 5.

2.6 Convergence Diagnostics

An important issue in applications of MCMC methods is to determine how many iterations of the algorithm are necessary before the Markov chain converges to its stationary distribution $\pi(\mathbf{x})$. It is important to note that this may not happen in a finite number of iterations, in which case convergence diagnostics are necessary in

order to assess how close we get. In addition, a burn-in period, which is the number of samples which is discarded before we think we are close enough to the target distribution to obtain ‘high-quality’ samples, is used. Even if we do begin sampling from the target distribution, we still need convergence diagnostics to determine how quickly we are sampling the state space of $\pi(\mathbf{x})$. If this sampling is occurring at a slow rate, we will need many samples to adequately explore the target distribution.

An excellent comparative review of thirteen well-known convergence diagnostics can be found in [13]. The methods described therein differ in a number of ways, including whether they produce quantitative or graphical results and use one or multiple chains; furthermore, some apply only to a subset of the existing MCMC methods. Cowles and Carlin [13] emphasize that “it is not possible to say with certainty that a finite sample from an MCMC algorithm is representative of an underlying stationary distribution” and suggest using a combination of techniques to assess MCMC convergence. These techniques include applying several of the convergence diagnostics (e.g., Gelman and Rubin statistic [14]) to 3–5 parallel chains as well as investigating autocorrelations and cross-correlations. Here, we describe three of the most common convergence diagnostics: autocorrelation plots, the Gelman and Rubin Statistic, and scaled regeneration quantile (SRQ) plots.

Autocorrelation Plots

Autocorrelation plots are one of the simplest convergence diagnostics. They test how fast we explore the state space once we have reached the stationary distribution. Autocorrelation is a measure of how well matched a signal is to a time-shifted version of itself as a function of the amount of the offset. In our case, the signal is stationary and corresponds to the sequence of samples generated by the MCMC algorithm and we are interested in measuring how fast the autocorrelation decays. Specifically, for a sequence of N samples $\{\mathbf{X}_i; i = 1, 2, \dots, N\}$, the autocorrelation $R(k)$ is defined as follows:

$$R(k) = \frac{1}{(N - k)} \sum_{t=1}^{N-k} (\mathbf{X}_t - \boldsymbol{\mu})(\mathbf{X}_{t+k} - \boldsymbol{\mu}),$$

where k is the lag, or the amount that the signal is time-shifted, and μ is the sample mean. The autocorrelation is normalized so that its values range from -1 to 1 as a function of the lag k . Ideally, as lag increases, the autocorrelation values will drop off quickly (e.g., exponentially), suggesting that the chain is mixing well. High autocorrelation values for large k are a signal that the chain is mixing slowly (i.e., a small part of the state space was explored after many iterations) and perhaps a different proposal density should be used. However, as will be shown in the example in Section 2.8, there are cases in which an inherent high correlation exists among the sample values, and a diagnostic such as this one may be highly uninformative.

Gelman and Rubin Statistic

The method proposed by Gelman and Rubin [14] was developed to help assess the convergence of chains underlying the Gibbs sampler and is based on a normal approximation to the posterior distribution. It includes two steps, the first of which is to create an over-dispersed estimate of the target distribution. This estimate is then used to start several independent parallel Markov chains (i.e., multiple Gibbs samplers are initialized according to this distribution). The second step analyzes these chains to form an estimate of the distribution of each scalar target random variable of interest. Both of these steps are described in more detail below.

The general method for finding the over-dispersed estimate of the target distribution has three steps. The first step involves finding the modes of the target distribution by approximating it using a mixture of multivariate normals, which can be fit to the density using the Expectation Maximization (EM) algorithm. In the second step, we obtain samples from the over-dispersed distribution by drawing samples from this normal mixture and scaling each sample vector by a positive scalar random variable (see [14] for exact details). This produces a new distribution which is a mixture of multivariate Student's t distributions. The third step involves sharpening this mixture of distributions, which is done by “downweighing” regions of low probability density with respect to the target distribution.

In the second part of Gelman and Rubin's method, we first draw m independent

points from the over-dispersed distribution determined in the first part and subsequently use these as starting values for m independent Markov chains, which we each run in parallel for $2n$ iterations. We discard the first n iterations of each chain to reduce the effect of the starting distribution on the sample values.

Then, for each scalar parameter of interest, denoted here by x , we calculate B and W as described by the equations below:

$$\frac{B}{n} = \sum_{i=1}^m \frac{\bar{x}_i - \bar{x}}{m-1}, \quad W = \sum_{i=1}^m \frac{s_i^2}{m},$$

where \bar{x}_i is the mean of the samples produced by the i^{th} chain, s_i^2 is the variance of the samples produced by the i^{th} chain, and \bar{x} is the mean of the latter n samples from all m sequences (\bar{x} is also the estimate of the target mean $\hat{\mu}$). We can then use B and W to estimate the target variance as follows:

$$\hat{\sigma}^2 = \frac{(n-1)W + B}{n}.$$

Now, we relax the estimate of the target distribution by allowing some variability in the estimates of the target mean and variance. Our new estimate takes the form of a Student's t distribution with center μ , scale $\sqrt{\hat{V}} = \sqrt{\sigma^2 + \frac{B}{mn}}$, and degrees of freedom df (see [14] for the exact formula for df). To monitor the convergence of the Markov chain, we estimate the factor by which the scale of the Student's t distribution for x may be reduced if we continued sampling indefinitely. This scale reduction factor \hat{R} is defined as follows:

$$\sqrt{\hat{R}} = \sqrt{\left(\frac{n-1}{n} - \frac{(m+1)B}{mnW} \right) \frac{df}{df-2}}.$$

The idea is that if this scale factor is high then running the Markov chain for more iterations may improve our inference about the target distribution. Therefore, we continue sampling until the scale reduction factors for all quantities of interest go to 1. This convergence diagnostic is included in the MrBayes software package, which also provides readily available algorithms for Bayesian inference in phylogeny.

Scaled Regeneration Quantile (SRQ) Plots

Another graphical convergence diagnostic called the scaled regeneration quantile (SRQ) plot was developed by Mykland *et al.* [15] as a means of detecting poorly-mixing chains. If we define T_i as the time when the chain returns to a specific state, where i ranges from 1 to the total number of states m , the SRQ plot is a plot of T_i/T_m versus i/m . The way to interpret this plot is that large deviations from the 45° line indicate a chain that is poorly mixing. Additionally, the slope of the line that connects the plot at i/m and j/m corresponds to the ratio of the estimated probability of the specific state from the entire chain to the estimate based on the part of the chain between i and j . The SRQ plot is one of the convergence diagnostics used by Li *et al.* [2].

2.7 Practical Considerations

When working with MCMC algorithms there are several practical points worth considering, including how to determine how many Markov chains to run and how to choose starting values for these chains. These and other topics are discussed in more detail below.

2.7.1 One Long Chain vs. Several Shorter Chains

The MCMC algorithm described in Section 2.5 is a single chain, but we can run multiple chains in parallel in order to improve the algorithm's exploration of the state space of the target distribution $\pi(\mathbf{x})$. Unfortunately, the literature offers conflicting recommendations on how many Markov chains should be run and for how long. Gelfand and Smith [16] suggest using many short chains. Gelman and Rubin [14] recommend using several long chains. Finally, Geyer [17] suggests running one very long chain.

Proponents of running several shorter chains argue that a comparison of these chains can reveal non-convergence and hence imply that the chain needs to be run for

longer in order to converge to the target distribution $\pi(\mathbf{x})$. However, this argument only seems to motivate running a longer chain in the first place since it has a higher chance of converging (i.e., the single long chain may actually converge even though the shorter chains imply non-convergence). This point implies that running several shorter chains can only be good to preemptively diagnose non-convergence. Further, if we find that our chain is mixing slowly then we should consider other probability models and proposal distributions for generating the Markov chain underlying the MCMC algorithm.

Another reason to run several shorter chains might be to aid in choosing a starting distribution that is spread out enough so that, in running these chains, all the high probability regions of the state space of $\pi(\mathbf{x})$ are explored. However, when dealing with complex problems in high-dimensional spaces, the chance of picking a starting distribution close to the target distribution is very small. Additionally, some run several shorter chains in order to produce independent samples. Since ergodic averaging (i.e. the sample mean converging in probability to the actual mean) does not require the samples to be independent, we do not need to change our sampling method for this reason.

Those who advocate for a single long chain argue that this chain has the best chance of finding new modes and that comparing several shorter chains can never prove convergence. Both of these statements, though accurate, should not be taken to mean that one single long chain is the right option. Several shorter runs can be used to determine how long we should run the single long chain. Once this information is known, we can run that single chain on the best probability model and use it to compute sample averages.

2.7.2 Optimizing the Acceptance Rate

One of the parameters we can control in the design of a Markov chain Monte Carlo algorithm is the acceptance rate, or the rate at which new samples proposed by our algorithm are kept. The acceptance rate is intimately related to the mixing rate described in Section 2.7.1. In general, a proposal distribution which generates

small steps $|X_i - Y|$ will have a high acceptance rate, or rate at which new samples are accepted, because new samples which are proposed are likely to remain in high probability regions. However, this same distribution may have a slow mixing rate, or rate at which the state space of the target distribution is explored, since the new samples are not likely to reach other modes of the distribution. On the other hand, a proposal distribution which generates large moves may have a low acceptance rate as the proposed moves are often ones from high to low probability regions of the target distribution. Ideally, we would like to avoid both scenarios and be able to design our proposal distribution in a way that optimizes the acceptance rate. Unfortunately, acceptance rates recommended in the literatures are only heuristics based on asymptotic behavior and simplifying assumptions. One such rule, developed by Roberts *et al.* [18], suggests using a proposal distribution with acceptance rate close to 0.25 for high dimensional models and 0.50 for models of dimension 1 or 2.

Overall, when choosing an acceptance rate, the higher level issue to consider is the mixing rate of the distribution. If the resulting behavior is undesirable, we should consider modifying the parameters of the proposal distribution or changing the proposal distribution to a different family of distributions. But generally, acceptance rates should be between 0.3 and 0.5.

2.7.3 Burn-in and Choosing a Starting Value

As described in Section 2.6, the burn-in period refers to the throwing away of a number of samples from the beginning of an MCMC run. The name burn-in comes from electronics. Since many electronic components tend to fail within a short span of time, a burn-in is done at the factory which aims to eliminate the worst components. Unfortunately, there is not a justification or motivation for burn-in in MCMC theory. The Markov chain versions of the Strong Law of Large Numbers and central limit theorem hold regardless of the starting distribution (see [3], Section 6.7). Burn-in, therefore, is simply a way of choosing a particular form for the starting distribution. In particular, if we start at a point whose probability distribution is π_0 and consider a length n burn-in period, then the starting distribution takes the following form:

$\pi_0 \mathbf{P}^n$, where \mathbf{P} is the transition probability matrix.

An obvious follow-up question is how to choose a starting value. One good rule of thumb well-summarized by Geyer is to “choose the starting point as any point you do not mind having in a sample.” So, ideally, we would choose points from higher probability regions of our target distribution. But, if we do not know anything about our target distribution, any point we choose is as good as any other.

2.7.4 Convergence Acceleration

There has also been work in developing methods to accelerate MCMC convergence. These methods include utilizing auxiliary variables [19,20], resampling and adaptive switching [21], simulated tempering [22], and Metropolis-coupled MCMC [23]. We describe Metropolis-coupled MCMC (also known as parallel tempering or swapping) below as it has often been used in applications to phylogenetics.

The motivation for using Metropolis-coupled MCMC (MC³) [23] and, in general, any convergence acceleration technique, comes from the fact that the target distribution $\pi(\mathbf{x})$ may have many modes separated by low probability regions, which the Markov chain underlying the MCMC algorithm may have difficulty traversing. In other words, the Markov chain can easily get stuck in one of the modes and the resulting samples will not accurately approximate the distribution of interest. The MC³ approach to this problem is to run m parallel Markov chains, each with a slightly different version of the target distribution, given by the following formula:

$$\pi_j(\mathbf{x}) = \pi(\mathbf{x})^{1/(1+\lambda(j-1))}, \quad (2.13)$$

where $\lambda > 0$ and $j \in \{1, 2, \dots, m\}$. The first chain $\pi_1(\mathbf{x})$ is the target density $\pi(\mathbf{x})$ and is referred to as the cold chain. The other chains are successively flatter versions of the target density (to see this, note that raising $\pi(\mathbf{x})$ to $\frac{1}{C}$ with $C > 1$ will flatten the distribution) and are referred to as heated chains¹. Since the heated chains

¹The naming convention of cold and heated chains comes from thermodynamics, where heating a metal has the effect of flattening it out.

correspond to progressively flatter versions of the target distribution, the underlying Markov chains have an increasingly easier time exploring the modes of the target distribution since proposed moves to originally low probability regions are now to higher probability regions, and hence are more likely to be accepted. After each iteration of the algorithm, we propose a swap between two randomly chosen chains i and k using a Metropolis-style step. Let \mathbf{X}_n^i be the state of chain i at iteration n and \mathbf{X}_n^k be the state of chain k at iteration n , where $i, k \in \{1, 2, \dots, m\}$ and $i \neq j$. Then, at iteration n , a state swap between chains i and k is accepted with probability:

$$\alpha(\mathbf{X}_n^i, \mathbf{X}_n^k) = \frac{\pi_i(\mathbf{X}_n^k)\pi_k(\mathbf{X}_n^i)}{\pi_i(\mathbf{X}_n^i)\pi_k(\mathbf{X}_n^k)}.$$

The pseudocode for the algorithm is shown in Algorithm 6.

At the end of the run, only the output of the cold chain is kept for analysis (i.e., the other chains are discarded). Since the hot chains will visit low probability regions of $\pi(\mathbf{x})$ more frequently than the cold chain, swapping states between chains will allow the cold chain to traverse the low probability regions of the target distribution more easily, which will lead to better overall mixing. However, as noted in [23], if the ratio $\frac{\pi_i(\mathbf{x})}{\pi_k(\mathbf{x})}$ is unstable (i.e., varies greatly depending on the particular chains which are chosen), proposed swaps between chains will not be accepted often. Hence, when implementing MC³, several hot chains that only differ by small increments should be used. Also, since a total of m chains are run, but only 1 (the cold one) is kept for analysis, this method has a clear computational disadvantage. Hence, MC³ implementations are best suited for machines with parallel computing capabilities.

2.8 Example: Sampling from a Bivariate Normal Distribution

In this section, we synthesize the preceding materials through an example, namely we show how the Metropolis-Hastings algorithm can be used to draw samples from a bivariate normal distribution. We draw samples using two different random walk

Algorithm 6: Metropolis-coupled MCMC

Input: $\pi(\mathbf{x})$

Result: N samples $\{\mathbf{X}_i; i = 1, 2, \dots, N\}$ that are distributed according to $\pi(\mathbf{x})$ as $N \rightarrow \infty$

begin

 Initialize m parallel chains according to their respective state spaces as described by Equation (2.13)

for $n = 1$ to N **do**

 Sample each of the m chains according to its target distribution $\pi_j(\mathbf{x})$

 Randomly choose two chains i and k

 Swap states of chains i and k with probability $\alpha(\mathbf{X}_n^i, \mathbf{X}_n^k)$

end

proposal densities, compare their results to those of a direct sampling method, and show autocorrelation plots in order to assess convergence. This example is based on the work of Chib and Greenberg [24]. We refer the reader to their paper for additional details, other potential proposal densities, and another example application of the Metropolis-Hastings algorithm.

We consider the problem of drawing samples from a bivariate normal distribution with mean vector $\boldsymbol{\mu} = (1, 2)^T$ and covariance matrix $\boldsymbol{\Sigma}$ shown below:

$$\boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}.$$

We can draw samples $\{\mathbf{X}_i; i = 1, 2, \dots, N\}$ from this distribution directly using the Cholesky decomposition as described below:

1. Determine \mathbf{P} such that it satisfies $\mathbf{P}^T \mathbf{P} = \boldsymbol{\Sigma}$ using Cholesky factorization.
2. Set $\mathbf{X}_i = \boldsymbol{\mu} + \mathbf{P}^T \boldsymbol{\tau}$, where $\boldsymbol{\tau}$ is a random variable distributed according to the bivariate normal distribution with mean vector $(0, 0)^T$ and identity covariance matrix.

This approach will be used to evaluate the results of the Metropolis-Hastings algorithm. We will generate samples from the bivariate normal distribution using two different symmetric random walk proposal densities of the form $\mathbf{Y}_i = \mathbf{X}_i + \mathbf{Z}$, listed below:

1. \mathbf{Z} is uniformly distributed, with the first component of \mathbf{Z} distributed uniformly on $[-0.75, 0.75]$ and the second component distributed uniformly on $[-1, 1]$.
2. \mathbf{Z} is distributed as an independent normal with mean vector $(0, 0)^T$ and a diagonal covariance matrix with the first diagonal entry 0.6 and the second diagonal entry 0.4.

For both of these proposal densities, the probability of a move is given by the following formula:

$$\alpha(\mathbf{X}_i, \mathbf{Y}_i) = \min \left(1, \frac{e^{-0.5(\mathbf{Y}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{Y}_i - \boldsymbol{\mu})}}{e^{-0.5(\mathbf{X}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{X}_i - \boldsymbol{\mu})}} \right). \quad (2.14)$$

We then generate 4,000 samples using the Cholesky approach and 6,000 samples using Metropolis-Hastings for each of the proposal densities. More observations are taken from the Metropolis-Hastings approach, as was done in [24] so that the scatter plots of the two approaches are comparable. Since our distribution has an inherently high correlation, as can be seen by the values on the diagonal of the covariance matrix, we expect these plots to show high autocorrelation, and hence will also consider other criteria (including the scatter plots previously mentioned) to assess the results of the Metropolis-Hastings algorithm. These criteria include sample mean and covariance, and 5th and 95th percentiles for the first and second dimension of the mean (using the first 4,000 samples for each approach). The k^{th} percentile number tells us that, for the particular dimension, $k\%$ of the samples had a value in that dimension below the specified number.

The scatter and autocorrelation plots are shown in Figure 2-2. As can be seen, for both proposal densities, the Metropolis-Hastings algorithm thoroughly explored the state space of the distribution in question. The autocorrelation plots, however, show that the sample correlation is quite high for both dimensions even for a lag of 600. Looking at these plots alone may lead us to believe that the chains are mixing slowly and therefore, we should draw more samples to ensure convergence. In this sense, the example does a good job of showing the dangers of relying too heavily on convergence diagnostics, which may falsely imply non-convergence.

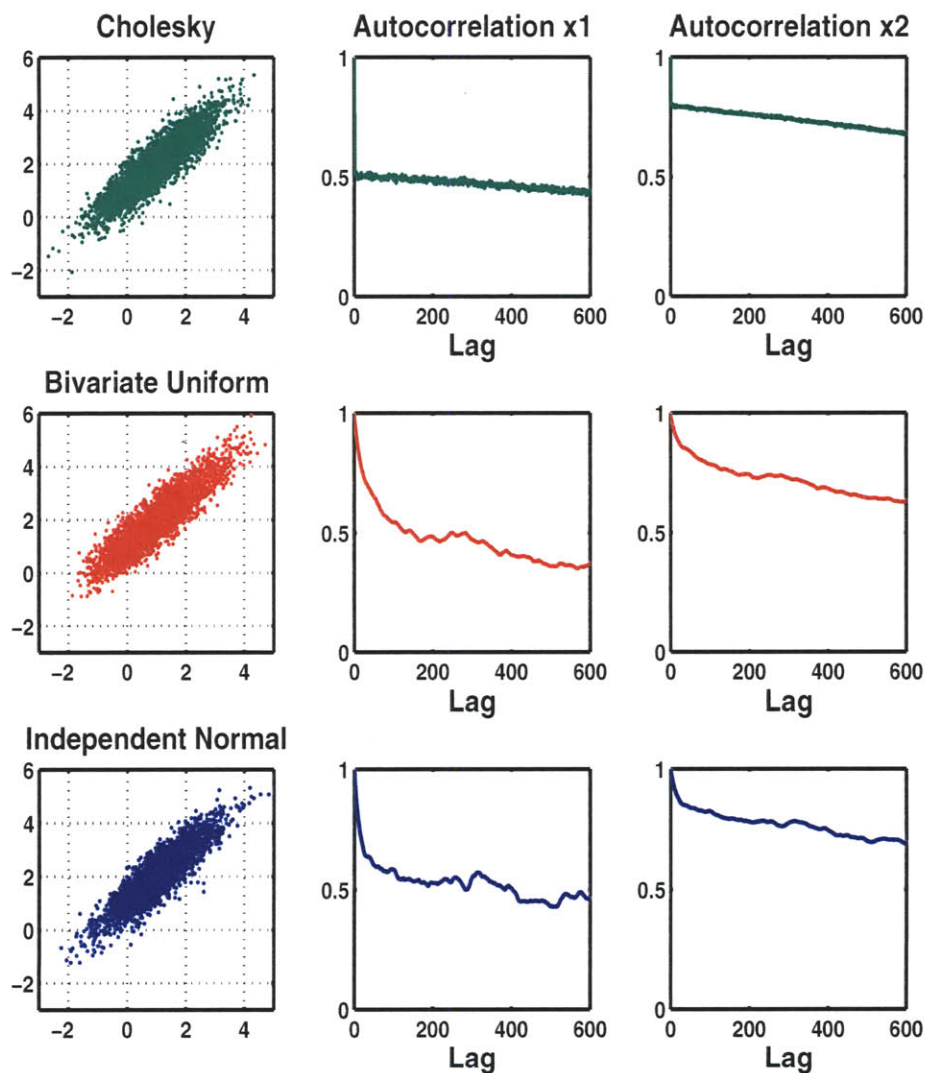


Figure 2-2: Diagnostic plots for simple MCMC algorithm for sampling from a bivariate normal distribution. The first column of figures shows 4,000 samples drawn using the direct sampling Cholesky approach and 6,000 samples drawn using the bivariate normal and independent normal proposal densities in the Metropolis-Hastings algorithm. The next two columns of figures show the autocorrelation plots as a function of lag for the first and second dimensions of the samples. As can be seen from the direct sampling approach, the large correlation is inherent in the data and is not actually a sign of non-convergence.

Method	<i>Cholesky</i>	<i>Bivariate Uniform</i>	<i>Independent Normal</i>
Mean	1.0081 2.0015	1.0807 2.0424	1.0772 2.0965
Covariance	0.9793 0.8856 0.8856 1.0021	1.0907 0.9806 0.9806 1.0842	0.9462 0.8558 0.8558 0.9556
Percentiles x_1	-0.6221 and 2.6472	-0.6080 and 2.9544	-0.4557 and 2.7186
Percentiles x_2	0.3958 and 3.6642	0.3303 and 3.7990	0.3303 and 3.7386

Table 2.1: Results of MCMC simulations from a bivariate Gaussian. The sample mean, covariance, and percentiles are shown for the three sampling methods. The two percentiles included for the first dimension x_1 and second dimension x_2 are those for the 5th and the 95th percentile. As can be seen from these values, the bivariate uniform and independent normal proposal densities produce comparable results.

Table 2.1 summarizes the results for the sample means, covariances, and percentiles computed for each sampling method. These values show that the Metropolis-Hastings algorithm has converged in both cases and produced quite accurate results, a conclusion which is supported by the scatter plots. It is further worth noting that we did take 4,000 samples from a distribution in a 2-dimensional space. The number of samples required for similar results in a higher-dimensional distribution underlying a more complex problem can be substantially higher. Also, as can be seen from the values for the percentiles, the bivariate uniform and independent normal proposal densities both produce similar results, hence this example does not motivate using one type of proposal density over the other.

Chapter 3

Phylogenetics

This chapter provides an introduction to phylogenetics, a field concerned with determining the evolutionary relationships among a group of organisms. Phylogeny refers to the origin and evolution of a set of organisms, in our case, a set of species. This evolution is often represented by a phylogenetic tree, which is described later in this chapter. We begin by describing the structure and evolution of genes. We proceed to discuss the role that phylogenetics plays in computational biology. In particular, we focus on phylogenetic tree construction. We first describe how phylogenetic trees relate to an underlying model for gene evolution. Then, we provide an overview of the common approaches to phylogenetic tree construction, including both distance and model-based approaches. Finally, we conclude with a discussion of various tree proposal methods, which enable us to explore the possible space of phylogenetic trees.

3.1 Genes and their Evolution

A gene is a hereditary unit consisting of a nucleic acid (e.g., DNA and RNA) sequence and determines a particular trait or characteristic of an organism. The nucleic acid is composed of nucleotides, which in the case of DNA are adenine, cytosine, guanine, and thymine, denoted by A, C, G, and T respectively. Adenine and guanine are called **purines** while thymine and cytosine are called **pyrimidines**. The genome of an organism refers to all of its hereditary information, which is encoded in the entirety

of its nucleic acid sequences (DNA for the case of most organisms).

New genes arise by two distinct processes. The first is called **speciation**, in which a species population diverges into multiple distinct species. In this case, each new species inherits a copy of the ancestral species' genes. Pairs of genes from these new genomes are referred to as **orthologs**. New genes can also arise through **gene duplication**, where a particular gene's nucleic acid sequence is copied and reinserted at another location within the genome. This process produces two copies of a gene within the same genome, which are called **paralogs**.

Biologists are interested in determining genes that are orthologs, since these genes often retain the same function after speciation. This means that if we know the function of a gene in one species, we can predict the function of another gene in a different species if the two genes in question are orthologs. However, we need to be careful since our prediction may be inaccurate in the presence of gene duplication. In this case, there are two copies of the same gene in the genome and one of the copies may evolve and adapt to a new function. Hence, it is important to identify both gene orthologs and paralogs for a set of species in order to improve the accuracy of our predictions. In the rest of this chapter, we will describe how we can perform this identification.

There has been a great increase in the amount of biological data available in recent years. Just in the last decade, we have sequenced a number of genomes, most recently including several species of flies and yeast, and more genomes are in the process of being sequenced. By obtaining this vast quantity of data, we have greatly increased our potential for understanding various aspects of the biology of cells and organisms. However, we face the challenge of how to extract meaningful information from this large amount of data. As noted by Felsenstein [25], the most reliable way to determine the structure and function of a biological molecule is through direct experimentation. However, it is far easier and less time consuming to obtain the DNA sequence of a particular gene. This fact motivates the development of computational methods to help infer biological information from sequence data.

One of the main inferential goals of these methods is to determine the evolutionary

relationships among a set of species, which then allows us to perform tasks such as determining the function of new genes. As described in [26], this process typically involves the following steps:

1. Annotate gene locations for a set of species
2. Reconstruct species phylogeny. Based on the genetic sequences from the previous step, a phylogenetic tree is constructed which describes the evolutionary relationships among the set of species in question.
3. Cluster gene families¹
4. Align gene families
5. Reconstruct gene phylogenies. For a set of modern genes, using their nucleotide sequences, reconstruct the phylogenetic tree which describes their evolutionary relationships.
6. Reconcile the gene phylogenies to the species phylogeny. For every gene phylogeny which was reconstructed, establish the mapping of genes in the gene tree to species in the species tree (i.e., determine the species corresponding to every gene).
7. Generate orthology database

An important point to clarify is the relationship between the species phylogeny in Step 2 and the gene phylogeny in Step 5. Both of these steps reconstruct phylogenetic trees and hence require the use of phylogenetic tree construction algorithms. The species phylogeny describes the evolutionary relationships among the set of species in question, and this, in turn, affects the structure of each gene phylogeny, which can be thought of as evolving within the species tree. In other words, the species phylogeny describes the general structure of how a gene evolves over time and macro-evolutionary events, such as gene duplications, determine the detailed structure of

¹A gene family is a set of genes that are descendants of a single gene in the last common ancestor (LCA), or most recent common ancestor w.r.t. evolutionary time, of all of the species in consideration

the tree. Step 6 uses these two ideas to determine the species of each gene in the gene phylogeny and hence, how the gene has evolved over time within the species phylogeny. Chapter 4 describes a model and phylogenetic tree construction algorithm that focuses on reconstructing species phylogenies while the model and algorithm of Chapter 5 focuses on reconstruction of gene phylogenies.

Collectively, this sequence of steps for the determination of evolutionary relationships among a set of species has been coined phylogenomics [27]. Here we focus on methods for reconstructing gene phylogenies, which includes performing Steps 5 and 6 above. In addition to phylogenetic tree construction, we also describe reconciliation, the process through which we relate each gene phylogeny to the species phylogeny, as it comprises an important part of our research.

3.2 Phylogenetic Tree Construction

A phylogenetic tree, which is an acyclic, connected graph of nodes, represents the evolutionary relationships within a group of organisms. Among many other applications, these phylogenetic trees can be used to reconstruct ancient proteins, study the function and properties of genes, aid in vaccine preparation, identify mutations associated with a particular disease, and determine the identity of new pathogens [1]. Nodes correspond to taxonomic units, such as species, populations, individuals, or genes. The nodes can either be **external** – the leaves of the tree – or **internal** – all the other nodes including the root node (if it exists). Associated with each node in the phylogenetic tree is a nucleotide sequence. The topology, or branching pattern of the tree, describes the relationships among these taxonomic units.

The phylogenetic trees that are most often considered in these studies are **bifurcating** trees in which each node has two children. A tree in which a node has more than two children is referred to as **multifurcating**. A phylogenetic tree may either be rooted or unrooted. In a **rooted** tree, there is a single node called the root which is the common ancestor of all of the nodes in the tree. Here, the flow of evolutionary time is along the path from the root node to all of the other nodes. In an **unrooted**

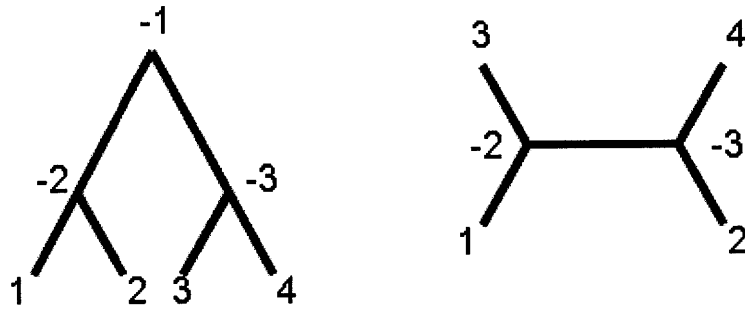


Figure 3-1: Examples of a rooted and an unrooted bifurcating phylogenetic tree. The tree on the left is rooted with -1 representing the root. The internal nodes are denoted by -2 and -3 while nodes $1 - 4$ are external nodes. The tree on the right is the same as the tree on the left except it has no root node, i.e., nodes $1 - 4$ are external nodes while -2 and -3 are internal nodes, as before. Note: nucleotide sequences at nodes are not shown for the sake of clarity.

tree, we only know the evolutionary relationships among the data without the knowledge of the evolutionary paths relating the nodes. Sometimes, an unrooted set of nodes is rooted by adding in a node, known as the **outgroup**, which is assumed to be a common ancestor to all the nodes under consideration. Examples of a rooted and an unrooted phylogenetic tree are shown in Figure 3-1.

The length of each branch (or edge between a parent and a child node) typically represents either the time over which the mutation took place (and can then be associated with a continuous-valued variable representing evolutionary time) or a number corresponding to the total number of mutations that took place in the transition from parent to child node². These mutations include observed mutations as well as **back mutations**, which are changes in a nucleotide pair that restore the original sequence. In this case, once the entire tree is constructed, it can be thought of as representing all the mutations which took place in order to produce the set of modern gene sequences (i.e., the external nodes of the tree) of interest.

Phylogenetic tree construction can be described as the search for the best explanatory model (tree) for an observed set of data (the external nodes of the tree, which, in our case, represent nucleotide sequences corresponding to modern genes). The se-

²We will see both types of trees in the following chapters and a rigorous mathematical formulation will be given at that time.

lection of the best explanatory tree is largely influenced by the underlying nucleotide sequence evolutionary model, which specifies how genetic sequences mutate over time. We provide an overview of these evolutionary models in the next section. The problem of phylogenetic tree construction is difficult in that there exists a large set of phylogenetic trees that are *consistent* with a particular set of modern genes. If we have $n > 1$ modern genes there are a total of $|\mathcal{T}| = \prod_{i=1}^{n-1} (2i - 1)$ possible distinctly labeled tree topologies (for proof, see Appendix A). If the sequences at each node contain m nucleotides and there are a total of 4 possible nucleotides (i.e., A, C, G, and T), then the total number of possible trees is given by $4^{m(n-2)} \prod_{i=1}^{n-1} (2i - 1)$. There are typically tens of genes with nucleotide sequences containing hundreds of bases associated with a single tree, but even for small values of these parameters, e.g., $n = 10$ and $m = 1$, the number of possible trees is extremely large – 2, 258, 332, 876, 800. Therefore, enumerating all the possible trees for a given set of modern genes is computationally infeasible for real data and so is the task of finding the “best” tree according to a given measure of optimality. To address this issue, various algorithms have been developed for finding an optimal tree. We describe them in greater depth in the next chapter. The focus will turn to using MCMC methods to design such algorithms.

3.2.1 Evolutionary Models

A crucial aspect of a phylogenetic tree is the underlying sequence evolution model describing how nucleotides may change over time. These models assume that nucleotides mutate stochastically and therefore attempt to model substitutions probabilistically (e.g., each nucleotide in the sequence has a certain probability of changing to one of the other three possible nucleotides and evolves independently of the other nucleotides in the same sequence). Another common assumption in these models is that each pair of branches which originate from a common ancestor evolve independently over time, which further reduces mathematical complexity.

The simplest model of DNA evolution was developed by Jukes and Cantor [28]. This is a one-parameter model in which nucleotide substitutions follow a Poisson process and each nucleotide can mutate into one of the other three with equal prob-

ability. Kimura [29] extended the model so that nucleotide changes still follow a Poisson process, but transitional events, or mutations within purines or pyrimidines have a different probability of occurrence from transversional events, or mutations between purines and pyrimidines.

The next extension came from Felsenstein [30]. He added three parameters to the original Jukes and Cantor model so that the user could specify a non-uniform distribution for the four nucleotide bases. A more recent generalization by Hasegawa *et al.* [31] combines the two previously mentioned extensions. The resulting model has five parameters and allows the distribution of the nucleotide bases to be non-uniform as well as the probabilities of transitional and transversional events to differ.

There are many other DNA sequence evolution models. For a more complete list, please see [25], Chapter 13. Chapter 14 of the same book describes protein evolution models, which are a counterpart to the DNA evolutionary models when a phylogenetic tree is used to model the evolution of amino acid sequences.

3.2.2 Reconciliation

Reconciliation provides a mapping from every gene in a gene phylogeny to some species in the species phylogeny. Typically the species of the modern genes are known while the species of the ancestral genes must be determined. If the genes only undergo sequence mutations, then the reconciliation of each gene is easily seen to be the last common ancestor (LCA), or the most recent common ancestor with respect to evolutionary time of all species under consideration, of the reconciliation of all of the children of that gene. Since we know the species for modern genes (and hence their reconciliation), we can determine the species of all of the parents of the modern genes by computing the LCA, and then repeat the procedure for the parents of these nodes and their parents until we have determined the species for every internal node in the tree.

However, as the species evolve, the genes can undergo macro-evolutionary events such as gene loss and duplication, in which case the simple LCA procedure has to be modified to account for these types of occurrences. The algorithm which determines

reconciliations under these circumstances often includes a penalty for each gene loss or duplication that is inferred (since reconciliations that infer large number of gene duplications and losses are less likely to be correct). Page [32] developed the first algorithm for finding the reconciliation which infers the fewest gene duplications. We rely on his algorithm to determine reconciliations in our work. While in Chapter 4 we discuss a model that does not deal with reconciliation, the information is incorporated into the model described in Chapter 5.

3.2.3 Distance-based Approaches

Here we survey a number of methods for phylogenetic tree estimation from data. We begin with distance-based methods, which rely on calculating the distance between each pair of modern species' sequences and then finding a tree whose structure most closely matches the observed distances. It is important to note that these approaches do not keep track of the specific types of mutations that occur (e.g., how many times substitutions there are from one base to another). We now describe three of the most popular algorithms in this category: neighbor-joining, minimum evolution, and parsimony. The latter two methods involve using an optimality criterion to search for the “best” tree. For a more complete list of distance-based approaches, we refer the reader to [1] and [25].

Neighbor-Joining In neighbor-joining, we first compute a matrix of observed distances \mathbf{O} for the n modern species, where entry \mathbf{O}_{ij} is the number of sites for which the nucleotides differ between the sequences of modern species i and modern species j . The algorithm, as described in [25], proceeds as follows:

1. For each tip (initially all of the modern species) u_i , compute $u_i = \sum_{j:j \neq i}^n \frac{\mathbf{O}_{ij}}{n-2}$.
2. Choose the i and j for which $\mathbf{O}_{ij} - u_i - u_j$ is the smallest. Join nodes i and j through a third node ij . Compute the branch lengths from node i to node ij

and from node j to node ij , respectively denoted as b_i and b_j , as follows:

$$b_i = \frac{1}{2}\mathbf{O}_{ij} + \frac{1}{2}(u_i - u_j), \quad b_j = \frac{1}{2}\mathbf{O}_{ij} + \frac{1}{2}(u_j - u_i).$$

3. Compute the distance between the new node ij and each of the remaining tips (indexed by k) as $\mathbf{O}_{ij,k} = \frac{\mathbf{O}_{ik} + \mathbf{O}_{jk} - \mathbf{O}_{ij}}{2}$.
4. Delete the entries for tips i and j from the matrix \mathbf{O} and replace them by an entry for the new node ij , which is now treated as a tip.
5. If more than two nodes remain, go back to Step 1. Otherwise, connect the two remaining nodes, say l and m by a branch of length \mathbf{O}_{lm} .

The tree produced by this method often becomes the input to an algorithm (e.g., parsimony) which tries to find the best phylogeny. The advantage of this method is ease of computation. As noted in [1], once the distance matrix is created, the neighbor-joining algorithm described in the steps above very quickly estimates a phylogenetic tree for the data. The disadvantage of this method, and other methods which use a distance matrix as input, is that the differences that are observed among sequences do not necessarily reflect the actual evolutionary distances³ between them. Sequences that have multiple substitutions at a site can have a very small observed and estimated evolutionary distance since a single substitution can often explain the differences in the two sequences. Therefore, the neighbor-joining method often estimates the incorrect tree when the gene sequences in question have large divergence times. On the other hand, when the sequences have diverged recently, neighbor-joining can be an extremely fast method for estimating a phylogenetic tree to describe the observed data.

Minimum Evolution The next two methods we describe are minimum evolution and parsimony, and both perform tree searches. These methods start with some initial

³The evolutionary distance is defined as the number of changes, including back mutations, that have occurred along the branches between two sequences.

estimate of a tree, which can be determined using an algorithm such as neighbor-joining, and then find the ‘best’ tree among a pre-defined set of trees relative to some optimality criterion uniquely defined for each method.

Minimum evolution is similar to neighbor-joining in that it requires a distance matrix representation of the data. The method seeks to create a tree consistent with the distance matrix that minimizes the total branch length of the tree (i.e., the sum of all of the branch lengths in the tree). This translates into finding what practitioners term the least squares tree for each possible topology considered, where the least squares measure is defined as follows:

$$Q = \sum_{i=1}^n \sum_{j=1}^n (\mathbf{O}_{ij} - o_{ij})^2,$$

where o_{ij} is the sum of the branch lengths between nodes i and j for a particular tree (and we are trying to minimize this quantity with respect to Q). As was shown in [25], we can rewrite this equation as follows:

$$Q = \sum_{i=1}^n \sum_{j:j \neq i} (\mathbf{O}_{ij} - \sum_k x_{ij,k} b_k)^2,$$

where $x_{ij,k}$ is an indicator variable equal to 1 if branch b_k lies on the path from species i to species j and 0 otherwise.

Now, we can differentiate Q with respect to each branch length b_k , set the derivatives to 0, and solve the resulting linear system to determine all the optimal branch lengths. The exploration of tree topologies can be performed by any of the methods described in the next section.

Similarly to neighbor-joining, the method is susceptible to failure because the estimated distance matrix does not account for multiple mutations at a single site. Recent studies have shown [1] that finding the optimal tree under the minimum evolution constraint does not provide a significant improvement over neighbor-joining, putting the utility of this method in question.

Parsimony Parsimony methods aim to quantify how well the mutations captured in the structure of the tree explain the resulting modern sequences. Each potential tree is assigned a score which represents the minimum number of mutations that could produce the modern sequences, and the tree with the lowest score is chosen. There are many versions of parsimony [25]. We will describe one of the simplest, called Wagner parsimony.

The Wagner parsimony method was developed for bifurcating trees with nucleotide sequence data, though it works for a general set of potential states of a site given that one state can change from any one to any other. The steps of the algorithm are as follows:

1. At each tip in the tree (initially all of the modern species), consider each group of sites⁴ one at a time. For each of these groups, create a set containing those nucleotides that are directly observed or are consistent with the previous observations. Thus, if we see an A at two sites we want to join (based on the topology), we create the set $\{A\}$. If we see an ambiguity, i.e., the two sites are different, we produce a new set which is the union of the two ambiguous sets.
2. Proceed up the tree (from the leaves to the root). At each internal node up to the root, determine the set that is the intersection of the sets of the two children nodes. If this set is empty, create a set that is the union of the sets at the child nodes and count one change of state.

Out of all the potential trees considered, the one with the smallest number of changes of state is chosen as the optimal tree. Once again, we can explore the possible space of tree topologies using any of the methods described in the next section. As outlined in [1], this method has several disadvantages. First of all, using the parsimony criterion leads to selecting a tree that has the fewest mutations among all of the paths from the ancestral sequences to each modern sequence. There are often many scenarios that could produce a specific set of modern sequences. A tree that has several, perhaps longer, pathways, most of which reasonably explain the presence

⁴A group of sites refers to the same site in each tip

of the modern sequences, should be chosen over a tree with only a few pathways, though shorter, that have the same explaining capability. Another disadvantage of this method is that no matter where a mutation occurs, the penalty added to the score is always the same. Therefore, two branches which are not adjacent on the actual tree can be evaluated as closest by parsimony. This is known as ‘long-branch attraction.’

3.2.4 Model-based Approaches

In recent years, model-based approaches to phylogenetic tree construction have gained popularity as they allow for the incorporation of sequence evolution models and, in the case of Bayesian models, allow for the quantification of the uncertainty associated with the underlying tree estimators. We describe the maximum likelihood and Bayesian model-based approaches below.

Maximum Likelihood The maximum likelihood method takes into account that multiple mutations can occur at a single site. Each proposed tree is evaluated based on how well it predicts the data (modern gene sequences), and the tree with the highest likelihood of producing the observed data is chosen. This approach relies on being able to calculate the conditional probability of seeing the data given a particular phylogenetic tree (i.e., the likelihood of the model). To do this, we must have a model, such as those described in Section 3.2.1, that describes the relative chance of various mutation events. This model can then be used to compute the probability of seeing the observed sequences given this particular tree.

As noted in [1], the advantage of the maximum likelihood method is that it considers all the possible topologies that are consistent with the data. Unfortunately, this fact makes the method computationally expensive, often too much so for practical applications when there is a large number of tree topologies consistent with the data. Also, as model complexity increases, finding the maximum likelihood tree becomes increasingly difficult. The search must be done on a multidimensional space of parameters and is not guaranteed to find the global optimum.

Bayesian Bayesian approaches to phylogenetic tree construction have become increasingly popular because they produce both a tree estimate and quantify the associated uncertainty. This allows us to compare various proposed trees, and produce a distribution on the set of possible trees as opposed to a point estimate whose relative likelihood is unknown. These approaches are closely related to the maximum likelihood methods and rely on Bayes' rule. The optimal tree hypothesis θ is the one that maximizes the hypothesis's posterior probability $P(\theta|Y)$, which is the probability of the hypothesis given the observed data Y (modern gene sequences). The prior probability $P(\theta)$ summarizes the belief that a hypothesis is the optimal one before having seen any data. Often, we do not know anything about the hypotheses a priori, so we choose an uninformative prior (e.g., uniform) and let the data inform the optimal solution.

As noted in Chapter 2, the difficulty in Bayesian methods typically lies in determining the normalizing constant in (2.1), since this requires computing the likelihood of every hypothesis in question. This time-intensive computation is avoided by using Markov chain Monte Carlo (MCMC) methods, described in Chapter 2, which can efficiently sample the posterior distribution without knowledge of the normalizing constant. Li *et al.* [2] used MCMC in just this way to develop an algorithm for phylogenetic tree construction.

An advantage of Bayesian approaches over traditional methods for phylogenetic tree construction is that they allow for the implementation of complex sequence evolution models which pose problems for the maximum likelihood approach. If the number of samples is small in comparison to the number of model parameters, maximum likelihood parameter estimates can be unreliable. In Bayesian analysis, the answer does not depend on one specific parameter value. Instead we marginalize over all of the possible parameter values. Even in the presence of adequate data, the hill-climbing algorithms that are used to find the optimal tree hypothesis in the maximum likelihood approach can be very slow and are not guaranteed to find the global optimum as the number of parameters increases [1].

3.2.5 Tree Proposal Methods

In many of the methods we described in the previous section, an important component to the underlying search is deciding how to propose new trees. Ideally, we would like to propose changes in tree topology that both allow us to quickly explore the entire space of possible tree topologies while at the same time proposing trees in frequency proportional to their likelihood, i.e., performing a probabilistic exploration of the space of trees. There are several common approaches to performing such an exploration.

Nearest-neighbor Interchange Nearest-neighbor interchange (NNI) is a search strategy which, in essence, swaps two adjacent branches of a tree in an attempt to propose a higher likelihood tree. The procedure is typically described for the case of an unrooted bifurcating tree. One version of this algorithm proceeds as follows and is shown in Figure 3-2:

1. Pick one of the internal branches (i.e., a branch connecting two internal nodes) of the tree. Erase the internal branch along with the two branches that are connected to it at each end (five branches in total).
2. The four resultant subtrees can be reconnected back into a tree in three different ways, one of which will produce the original tree. Try both of the other rearrangements and pick the tree which has the higher likelihood.
3. Repeat Steps 1 and 2 for every other internal branch until there is no further rearrangement that yields a higher likelihood tree.

Figure 3-2 shows the original tree and the two possible alternate rearrangements. If we have n external nodes, then there are $n - 3$ internal branches. Hence, in a single pass of the NNI algorithm $2(n - 3)$ trees are examined (since two different tree rearrangements are considered for each internal branch). Alternative versions of the algorithm simultaneously try rearrangements on multiple trees if they are tied for best in Step 2, the most sophisticated of which will keep track of all of the trees tied for best and rearrange all of them.

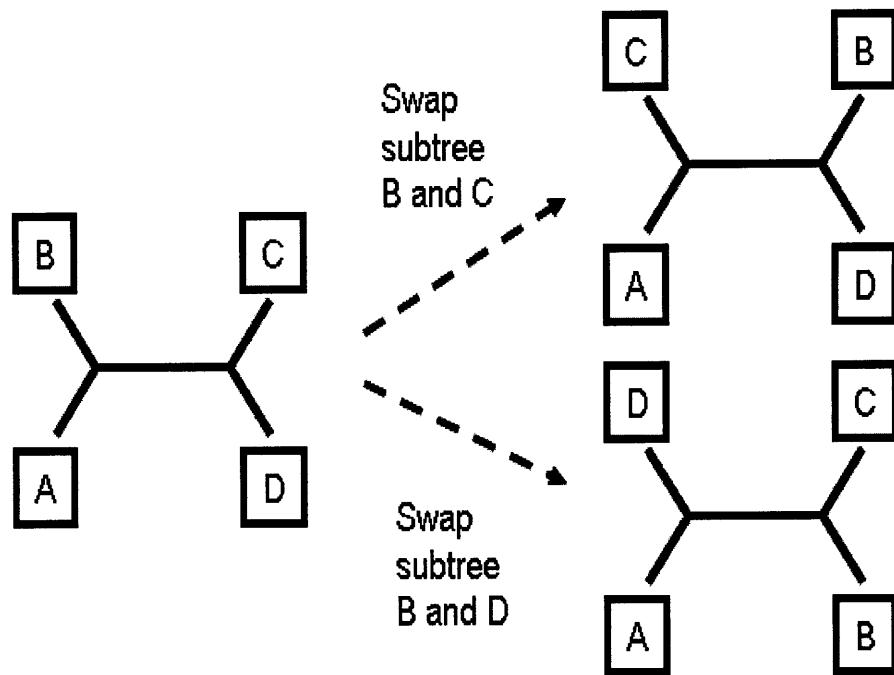


Figure 3-2: Swapping procedure in nearest-neighbor interchange algorithm. The original tree is on the left. The two trees on the right represent the two possible resulting subtrees post rearrangement.

Subtree Pruning and Regrafting Subtree pruning and regrafting (SPR), as the name suggests, involves removing an interior or exterior branch from the tree along with the attached subtree (pruning) and then reinserting the subtree back into the tree (regrafting). For an unrooted bifurcating tree with n external nodes, the total number of neighboring trees examined by SPR is $4(n - 2)(n - 3)$. Of course, some of these trees are the same as the original tree; nonetheless, SPR examines a much wider range of trees than NNI does. As in NNI, the method can simply accept the most likely tree for a particular pruned subtree or delay accepting a new tree until other SPR arrangements have been considered.

Tree Bisection and Reconnection Tree bisection and reconnection (TBR) is the most elaborate of the tree rearrangement strategies of those described so far. In this method, an interior branch is chosen and broken. The two resulting tree fragments are considered separate trees. Now, the method tries all possible connections between a branch from one tree and a branch from another tree. As in the other methods, one

of these connections will result in the original tree. The total number of trees that is considered depends on the particular topology of the tree (i.e., there is no closed-form formula). Once again, there is a choice in how greedy the method will be, i.e., it can simply accept the best of these re-connections or keep track of a list of the most likely re-connections as new rearrangements are tried.

There are many other possible tree proposal methods. The local rearrangement strategy of Li *et al.* [2] is described in the next chapter. Alternative options include tree-fusing, genetic algorithms, tree windows, and sectorial search. Descriptions of all of these algorithms can be found in [25].

Chapter 4

Phylogenetic Tree Construction Using Markov Chain Monte Carlo

The chapter discusses the work of Li *et al.* [2], who applied Markov chain Monte Carlo methods to phylogenetic tree construction. Their work is a good example of MCMC applied to phylogenetics and comprises the basis for the model in Chapter 5. We refer the reader to [33–35] for other MCMC approaches to phylogenetics. In the sections below, we describe their approach and our implementation of their algorithms well as provide a review of the results.

4.1 Motivation and Model Description

Li *et al.* in [2] were interested in the problem of how to reconstruct phylogenetic relationships between different organisms given their DNA or nucleotide sequences. Their Bayesian approach, described in detail below, includes an MCMC algorithm, which produces a distribution on the possible set of phylogenetic trees for a given set of nucleotide sequences.

4.1.1 Model Parameters

We denote the set of parameters characterizing the model of Li *et al.* by $\theta = \{\mathcal{D}, \mathcal{T}, \mathcal{E}\}$. Here \mathcal{T} corresponds to the set of possible distinct labeled topologies of bifurcating rooted trees¹. Each node in the tree represents a nucleotide sequence corresponding to a particular organism and the topology of the tree represents the evolutionary history of this set of organisms. If we have n taxonomic units (i.e., n nucleotide sequences whose evolutionary history we wish to recreate), then the total number of distinct labeled topologies $|\mathcal{T}|$, or branching patterns with n external nodes, of rooted bifurcating trees is given by the following formula²:

$$|\mathcal{T}| = \prod_{i=1}^{n-1} (2i - 1).$$

The nucleotide sequences for all of the nodes in the tree are denoted by \mathcal{D} . Let D denote the set of the four possible bases (i.e., $D = \{A, C, G, T\}$). We will denote each m -base nucleotide sequence by $\mathbf{d}_i \in D^m$ where $i \in \{-1, -2, \dots, -(n-1)\}$ indexes internal nodes (with -1 corresponding to the root) and $i \in \{1, 2, \dots, n\}$ indexes the external nodes. This means there are $4^{m(2n-1)}$ possible nucleotide sequences (since there are n external nodes and $n-1$ internal nodes each of which has a length m nucleotide sequence). Therefore, we can represent the set of all possible nucleotide sequences as $\mathcal{D} = (\mathbf{d}_{-1}, \mathbf{d}_{-2}, \dots, \mathbf{d}_{-(n-1)}, \mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_n)$.

In a rooted tree, a unique path exists from the root to each of the external nodes. The evolutionary time that passes from the root node to each of the external nodes (i.e., the total branch length) is assumed to be a constant $\tau > 0$. The convention is to assign the root an evolutionary time of τ and assign all of the external nodes an evolutionary time of zero. The evolutionary times of the remaining $n-2$ nodes are denoted by $t_i \in (0, \tau)$ where $i \in (-2, -3, \dots, -(n-1))$. Collectively, the evolutionary times for these nodes are denoted by $\mathcal{E} = (t_{-2}, t_{-3}, \dots, t_{-(n-1)})$.

The underlying nucleotide evolution model is based on the five-parameter model

¹All the phylogenetic trees Li *et al.* consider are rooted bifurcating trees, although they describe how to extend their methods to more complicated cases.

²For proof, see Appendix A

of Hasegawa *et al.* [31]. The overall transversion rate is represented by $\alpha \in \mathbb{R}_+$. The relative bias of transitional events over transversional events is represented by $\kappa \in \mathbb{R}_+$. The probabilities of each of the 4 possible nucleotide bases are represented by π_A , π_C , π_G , and π_T , which must be non-negative numbers that sum to 1. We will describe the process of nucleotide substitutions for this model using standard continuous-time discrete state space Markov chain theory. The nucleotide sequence evolution is governed by the generator matrix \mathbf{R} specified below:

$$\frac{\mathbf{R}}{\alpha} = \begin{pmatrix} -(\kappa\pi_G + \pi_C + \pi_T) & \pi_C & k\pi_G & \pi_T \\ \pi_A & -(\kappa\pi_T + \pi_A + \pi_G) & \pi_G & k\pi_T \\ k\pi_A & \pi_C & -(\kappa\pi_A + \pi_C + \pi_T) & \pi_T \\ \pi_A & k\pi_C & \pi_G & -(\kappa\pi_C + \pi_A + \pi_G) \end{pmatrix}, \quad (4.1)$$

where each entry \mathbf{R}_{ij} corresponds to the probability of a single site mutating from base i to base j , where bases A, C, G, and T correspond to rows and columns 1, 2, 3, 4, respectively.

The matrix of transition probabilities at time t , denoted by $Q(t)$ (where entry $Q_{ij}(t)$ represents the probability of transition from base i to base j in time t), for this Markov chain is given by the matrix exponential $e^{\mathbf{R}t}$. For the five-parameter model, the entries are given by the following formula:

$$Q_{ij}(t) = \begin{cases} \pi_j + \pi_j \left(\frac{1}{\lambda_j}\right) e^{-\alpha t} + \left(\frac{\lambda_j - \pi_j}{\lambda_j}\right) e^{-\alpha\gamma_j t} & i = j \\ \pi_j + \pi_j \left(\frac{1}{\lambda_j}\right) e^{-\alpha t} - \left(\frac{\pi_j}{\lambda_j}\right) e^{-\alpha\gamma_j t} & i \neq j \text{ (transitional event)} \\ \pi_j (1 - e^{-\alpha t}), & i \neq j \text{ (transversional event)} \end{cases}$$

where $\lambda_j = \pi_A + \pi_G$ if base j is a purine (A or G) and $\lambda_j = \pi_C + \pi_T$ if base j is a pyrimidine (C or T), and $\gamma_j = 1 + (\kappa - 1)\lambda_j$. In the Jukes-Cantor model, which was described in Section 3.2.1 and is used throughout this thesis, the equation above

simplifies as follows:

$$Q_{ij}(t) = \begin{cases} \frac{1}{4} + \frac{1}{2}e^{-\alpha t} + \frac{1}{2}e^{-\alpha t} & i = j \\ \frac{1}{4} + \frac{1}{2}e^{-\alpha t} - \frac{1}{2}e^{-\alpha t} & i \neq j \text{ (transitional event)} \\ \frac{1}{4}(1 - e^{-\alpha t}). & i \neq j \text{ (tranversional event)} \end{cases}$$

The model relies on the following assumptions:

- Given a topology \mathcal{T} , the evolutionary times of the internal nodes \mathcal{E} are uniformly distributed over the space of possible times, subject to the ordering constraints induced by the topology.
- The evolution of nucleotides at different sites and at different branches are independent random processes.

We can now construct a precise mathematical definition for a phylogenetic tree. As we have already mentioned, each tree is defined by the parameters \mathcal{T} , \mathcal{D} , and \mathcal{E} . Therefore, the probability distribution on the set of possible trees is the joint probability $P(\mathcal{D}, \mathcal{E}, \mathcal{T})$. We can rewrite this joint probability as the product of the following three terms:

$$P(\mathcal{D}, \mathcal{E}, \mathcal{T}) = P(\mathcal{D}|\mathcal{E}, \mathcal{T})P(\mathcal{E}|\mathcal{T})P(\mathcal{T}).$$

4.1.2 Model Specification

We now describe how to compute each of the three terms above in turn, proceeding from the right.

Prior distribution of topologies: $P(\mathcal{T})$. We assume that the distribution on topologies is uniform over all possible topologies, that is, each topology is equally likely a priori with the probability given by:

$$P(\mathcal{T}) = \frac{1}{\prod_{i=1}^{n-1} (2i - 1)}.$$

Conditional distribution of evolutionary times given topology: $P(\mathcal{E}|\mathcal{T})$.

Once we know the topology, our first assumption tells us that the evolutionary times are uniformly distributed over $(0, \tau)$ with the times of the parent nodes strictly larger than those of their offspring. If there is no parent/child relationship between two nodes, then their times are independently distributed given their most recent common ancestor. This results in the following distribution on times, for $n \geq 2$:

$$P(\mathcal{E}|\mathcal{T}) = \frac{1}{\tau^{n-2}} \prod_{i=2}^{n-1} (n_{-i} + 1),$$

where n_{-i} is the number of descendant internal nodes of internal node $-i$ for $i \in \{2, 3, \dots, n-1\}$.

Conditional distribution of nucleotide sequences given evolutionary times and topology: $P(\mathcal{D}|\mathcal{E}, \mathcal{T})$. The last term we need to compute is the likelihood of the nucleotide sequences given the evolutionary times and the topology. In order to compute this term, we need to be able to determine the probability $P_{\mathbf{d}_k \mathbf{d}_l}(t)$, which is the probability that the nucleotide sequence at node k changes to the nucleotide sequence at node l over a period of time t . There are 16 possible categories of nucleotide substitutions across the m sites in the nucleotide sequences at the two nodes. Let the number of sites with nucleotide i in sequence \mathbf{d}_k and nucleotide j in sequence \mathbf{d}_l be denoted by m_{ijkl} . Then the probability $P_{\mathbf{d}_k \mathbf{d}_l}(t)$ can be written as follows:

$$P_{\mathbf{d}_k \mathbf{d}_l}(t) = \prod_{i,j \in D} Q_{ij}(t)^{m_{ijkl}}. \quad (4.2)$$

This is a generic form for the probability of observing the nucleotide substitutions between the two sequences during a period of time t . Depending on the nucleotide evolution model (see Section 3.2.1 for a description of the common evolutionary models), the particular form of $Q_{ij}(t)$ may change.

Since we have made the assumption that along each branch of the tree the nucleotide sequences evolve independently, the probability of all the sequences in the tree, given the evolutionary times and topology, is proportional to the product of

the probabilities of all the nucleotide mutations from every parent to child pair in the tree. To complete the specification of the distribution of nucleotide sequences given evolutionary times and topology, we need to include the probability of starting at the root, which is a product of the probabilities $\pi_A, \pi_C, \pi_G, \pi_T$ of each site in the nucleotide sequence at the root. If we let r_j be the number of sites in the root node for which the nucleotide is j , where $j = A, C, G, \text{ or } T$, then the probability of starting at the root is $\prod_{j \in D} \pi_j^{r_j}$, and the overall likelihood that we are interested in is the following:

$$P(\mathcal{D}|\mathcal{E}, \mathcal{T}) = \left(\prod_{j \in D} \pi_j^{r_j} \right) \prod_{k, l \in E_{\mathcal{T}}} P_{\mathbf{d}_k \mathbf{d}_l}(t_k - t_l),$$

where $E_{\mathcal{T}}$ is the set of edges in \mathcal{T} and t_k and t_l are the evolutionary times at nodes k and l respectively.

4.2 Inference Algorithm

In this section, we describe Li *et al.*'s [2] algorithm for phylogenetic tree construction. We begin with the motivation for developing the algorithm and then proceed to the algorithm's description.

4.2.1 Motivation for Developing an MCMC-based Inference Algorithm

For a given problem, the observed variables are the nucleotide sequences for the external nodes and the root, which we will denote as $\mathcal{D}_{ext} = (\mathbf{d}_{-1}, \mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_n)$. The rest of the nucleotide sequences will be denoted by $\mathcal{D}_{int} = (\mathbf{d}_{-2}, \mathbf{d}_{-3}, \dots, \mathbf{d}_{-(n-1)})$ so that $\mathcal{D} = \mathcal{D}_{int} \cup \mathcal{D}_{ext}$. The target distribution $\boldsymbol{\pi}$ is the conditional distribution of the tree given the observed nucleotide sequences, i.e., $\boldsymbol{\pi} = \sum_{\mathbf{d}_{-1} \in D^m} P(\mathcal{T}, \mathcal{E}, \mathcal{D}_{int} | \mathcal{D}_{ext})$, and run an MCMC algorithm based on Metropolis-Hastings whose underlying Markov chain produces trees from this distribution of interest. The motivation for using MCMC comes from the fact that the normalizing constant for $\boldsymbol{\pi}$ is a sum of $\frac{1}{\prod_{i=1}^{n-1} (2i-1)}$ high-dimensional integrals.

4.2.2 Algorithm Description

Algorithm 7: Li *et al.*'s Algorithm

Data: \mathcal{D}_{ext}

Result: N trees $\{\mathbf{X}_i; i = 1, 2, \dots, N\}$ whose limiting distribution is π

begin

 Initialize \mathbf{X}_0 as described in Section 4.3

for $i = 1$ to N **do**

 Sample a tree \mathbf{Y}_i from $q(\mathbf{y}|\mathbf{X}_{i-1})$ as follows:

 1. Perform Li *et al.*'s local rearrangement strategy to produce a new tree \mathbf{Y}_i

 2. Determine a new time for the target node in \mathbf{Y}_i

 3. Determine a new sequence for target node in \mathbf{Y}_i

 Sample a uniform random variable U on $[0, 1]$

if $U \leq \alpha(\mathbf{X}_i, \mathbf{Y}_i)$ **then**

 | $\mathbf{X}_{i+1} = \mathbf{Y}_i$

else

 | $\mathbf{X}_{i+1} = \mathbf{X}_i$

end

We will now outline the algorithm Li *et al.* [2] developed for phylogenetic tree construction based on the model we have previously described. The pseudocode for the algorithm is shown in Algorithm 7 with a more detailed description provided below.

The algorithm employs a local rearrangement strategy in the neighborhood of one of the internal nodes T (called the target) excluding the root. This neighborhood includes the target's parent P , sibling S , and children $C1$ and $C2$. The algorithm based on the local rearrangement strategy can be described in the following steps:

- **Step 1:** Randomly choose a node among the $n - 2$ internal nodes (excluding the root) to be the target node T . From among $C1$, $C2$, and S , randomly select one of the nodes to be the new sibling S' . The remaining two nodes become the new children - $C1'$ and $C2'$.
- **Step 2:** Given the topology selected in Step 1, pick a time for the new target T' , which ranges between the time of the parent t_P and the more recent of the children's times, i.e. $\max\{t_{C1'}, t_{C2'}\}$. This time is picked from the following

density:

$$g(t_{T'}) = \frac{\sum_{\mathbf{d} \in D^m} P_{\mathbf{d}P\mathbf{d}}(t_P - t_{T'}) P_{\mathbf{d}\mathbf{d}_{C1'}}(t_{T'} - t_{C1'}) P_{\mathbf{d}\mathbf{d}_{C2'}}(t_{T'} - t_{C2'})}{\int_{\max\{t_{C1'}, t_{C2'}\}}^{t_P} \sum_{\mathbf{d} \in D^m} P_{\mathbf{d}P\mathbf{d}}(t_P - t) P_{\mathbf{d}\mathbf{d}_{C1'}}(t - t_{C1'}) P_{\mathbf{d}\mathbf{d}_{C2'}}(t - t_{C2'}) dt}$$

for $\max\{t_{C1'}, t_{C2'}\} \leq t_{T'} \leq t_P$

- **Step 3:** Having picked the time of the new target node, the last step of the algorithm is determining a new nucleotide sequence for the target node, which we will denote by $\mathbf{d}_{T'}$. The distribution on this sequence is as follows:

$$h(\mathbf{d}_{T'}) = \frac{P_{\mathbf{d}P\mathbf{d}_{T'}}(t_P - t_{T'}) P_{\mathbf{d}_{T'}\mathbf{d}_{C1'}}(t_{T'} - t_{C1'}) P_{\mathbf{d}_{T'}\mathbf{d}_{C2'}}(t_{T'} - t_{C2'})}{\sum_{\mathbf{d} \in D^m} P_{\mathbf{d}P\mathbf{d}}(t_P - t_{T'}) P_{\mathbf{d}\mathbf{d}_{C1'}}(t_{T'} - t_{C1'}) P_{\mathbf{d}\mathbf{d}_{C2'}}(t_{T'} - t_{C2'})}$$
 for $\mathbf{d}_{T'} \in D^m$

Now if we denote our initial tree by \mathbf{X}_i and the tree produced after the rearrangement procedure described above as \mathbf{Y}_i , then the Markov chain transition probability density for the local rearrangement strategy is $P(\mathbf{X}_i, \mathbf{Y}_i) = \alpha(\mathbf{X}_i, \mathbf{Y}_i)q(\mathbf{Y}_i|\mathbf{X}_i)$, with

$$\alpha(\mathbf{X}_i, \mathbf{Y}_i) = \min \left(\frac{\pi(\mathbf{Y}_i)q(\mathbf{X}_i|\mathbf{Y}_i)}{\pi(\mathbf{X}_i)q(\mathbf{Y}_i|\mathbf{X}_i)}, 1 \right), \quad (4.3)$$

$$q(\mathbf{Y}_i|\mathbf{X}_i) = P(T'|T)P(\mathcal{E}'|T', \mathcal{E})P(\mathcal{D}'_{int}|\mathcal{D}_{int}, T', \mathcal{E}'), \quad (4.4)$$

$$q(\mathbf{X}_i|\mathbf{Y}_i) = P(T|T')P(\mathcal{E}|T, \mathcal{E}')P(\mathcal{D}_{int}|\mathcal{D}'_{int}, T, \mathcal{E}), \quad (4.5)$$

where ' signifies a post local rearrangement parameter.

Based on the model description above and the properties satisfied by the underlying Markov chain (see [2] for details), the equations above become:

$$\alpha(\mathbf{X}_i, \mathbf{Y}_i) = \min \left(\frac{P_{\mathbf{d}P\mathbf{d}_{S'}}(t_P - T_{S'})(n_{T'} + 1) \int_{\max\{t_{C1'}, t_{C2'}\}}^{t_P} \sum_{\mathbf{d} \in D^m} P_{\mathbf{d}P\mathbf{d}}(t_P - t) P_{\mathbf{d}\mathbf{d}_{C1'}}(t - t_{C1'}) P_{\mathbf{d}\mathbf{d}_{C2'}}(t - t_{C2'}) dt}{P_{\mathbf{d}P\mathbf{d}_S}(t_P - T_S)(n_T + 1) \int_{\max\{t_{C1'}, t_{C2'}\}}^{t_P} \sum_{\mathbf{d} \in D^m} P_{\mathbf{d}P\mathbf{d}}(t_P - t) P_{\mathbf{d}\mathbf{d}_{C1'}}(t - t_{C1'}) P_{\mathbf{d}\mathbf{d}_{C2'}}(t - t_{C2'}) dt}, 1 \right)$$

$$q(\mathbf{Y}_i|\mathbf{X}_i) = \frac{P_{\mathbf{d}P\mathbf{d}_T}(t_P - t_T) P_{\mathbf{d}_T\mathbf{d}_{C1}}(t_T - t_{C1}) P_{\mathbf{d}_T\mathbf{d}_{C2}}(t_T - t_{C2})}{3^{(n-2)} \int_{\max\{t_{C1'}, t_{C2'}\}}^{t_P} \sum_{\mathbf{d} \in D^m} P_{\mathbf{d}P\mathbf{d}}(t_P - t) P_{\mathbf{d}\mathbf{d}_{C1}}(t - t_{C1}) P_{\mathbf{d}\mathbf{d}_{C2}}(t - t_{C2}) dt}$$

$$q(\mathbf{X}_i|\mathbf{Y}_i) = \frac{P_{\mathbf{d}P\mathbf{d}_{T'}}(t_P - t_{T'}) P_{\mathbf{d}_{T'}\mathbf{d}_{C1'}}(t_{T'} - t_{C1'}) P_{\mathbf{d}_{T'}\mathbf{d}_{C2'}}(t_{T'} - t_{C2'})}{3^{(n-2)} \int_{\max\{t_{C1'}, t_{C2'}\}}^{t_P} \sum_{\mathbf{d} \in D^m} P_{\mathbf{d}P\mathbf{d}}(t_P - t) P_{\mathbf{d}\mathbf{d}_{C1'}}(t - t_{C1'}) P_{\mathbf{d}\mathbf{d}_{C2'}}(t - t_{C2'}) dt}$$

The parameters underlying the generator matrix in Equation (4.1) are often unknown and must be estimated as part of the phylogenetic tree construction algorithm. Li *et al.*'s approach is to re-estimate these parameters during the burn-in period of their algorithm. The details of this procedure can be found in [2] along with a description of all of the convergence diagnostics that they implemented.

4.3 Basic Implementation and Results

We implemented the algorithm as described in Algorithm 7 with two exceptions: we held time constant, as this was more relevant to our approach in which branch lengths encapsulate a per site mutation rate, and we did not re-estimate parameters during the burn-in period of the algorithm. The parameters were held constant at the following values: all the π_i 's were 0.25, $\tau = 4$, $\alpha = 2$, $\kappa = 0.5$, $n = 4$, and $m = 2$.

We generated the initial tree for the algorithm using the generative model described in Section 4.1.2. We first generated a topology uniformly over the set of possible topologies by randomly choosing two subtrees (initially the subtrees are just the external nodes of the tree) to join and repeatedly joining subtrees until only two subtrees remain, which are joined by the root node. Once we had a topology, we fixed the times of the tree as follows: the root node's evolutionary time was set to τ and the times of all the external nodes were set to 0. The times of all of the nodes at level l below the root was set to $\tau - 0.1l$ with the exception of the external nodes, whose times remained unchanged. When we ran the algorithm, we updated the times accordingly to meet these constraints post local rearrangement. Finally, using the independent site evolution assumption, we generated the m -length sequences for all of the internal nodes of the tree, one site at a time, starting from the children of the root node and proceeding down the tree, using Equation 4.2³.

We ran 20 parallel Markov chains according to algorithm described in Algorithm 7 for 10,000 samples each and computed errors averaged over all of the chains to ensure

³For each possible nucleotide base, we compute the probability that the parent base changes to this base over the time that elapsed and then we choose the new base according to the resulting discrete probability distribution.

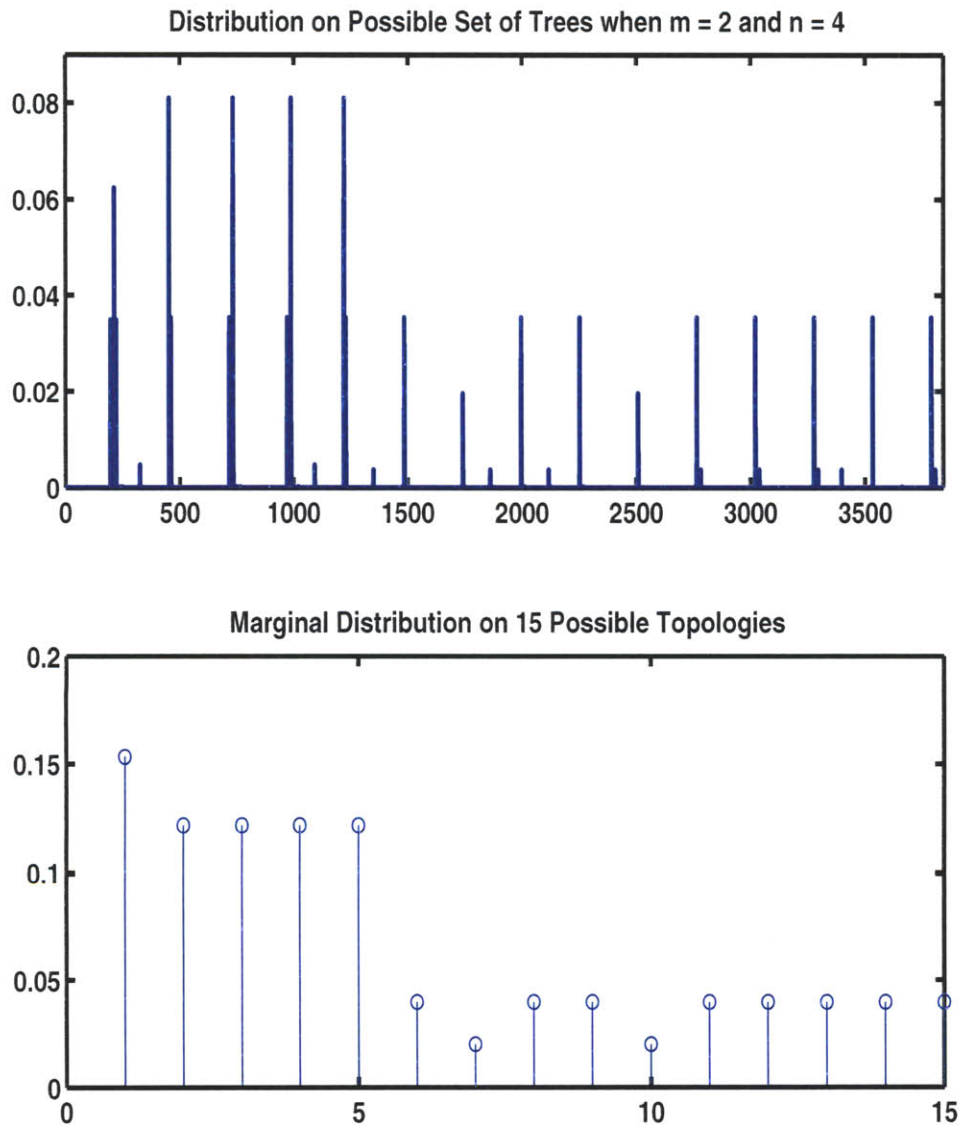


Figure 4-1: Distribution on the set of possible topologies. Since $n = 4$ and $m = 2$, there are 3,840 possible trees. Their relative probabilities are shown in the top graph. The bottom graph shows the marginal probabilities of the 15 possible topologies.

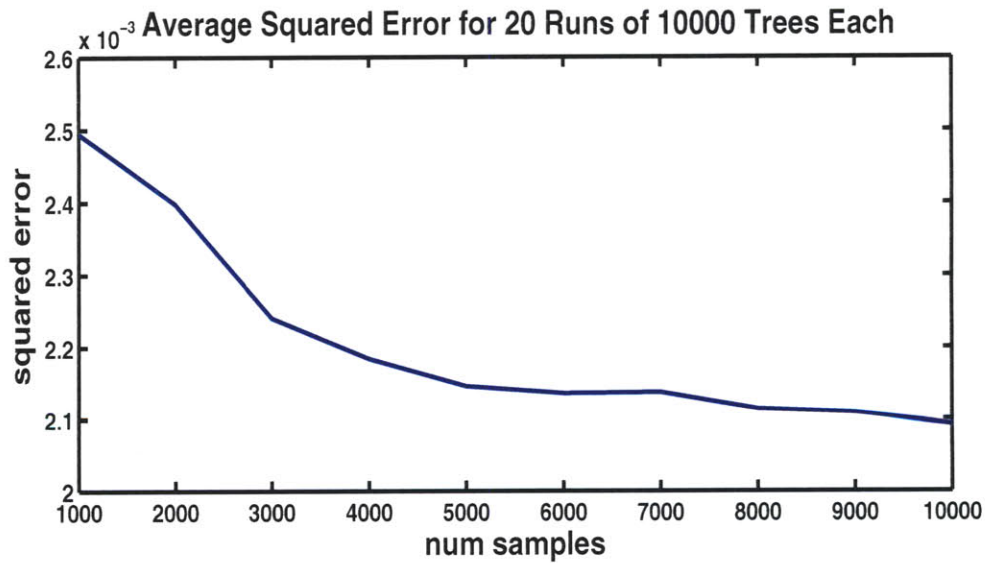
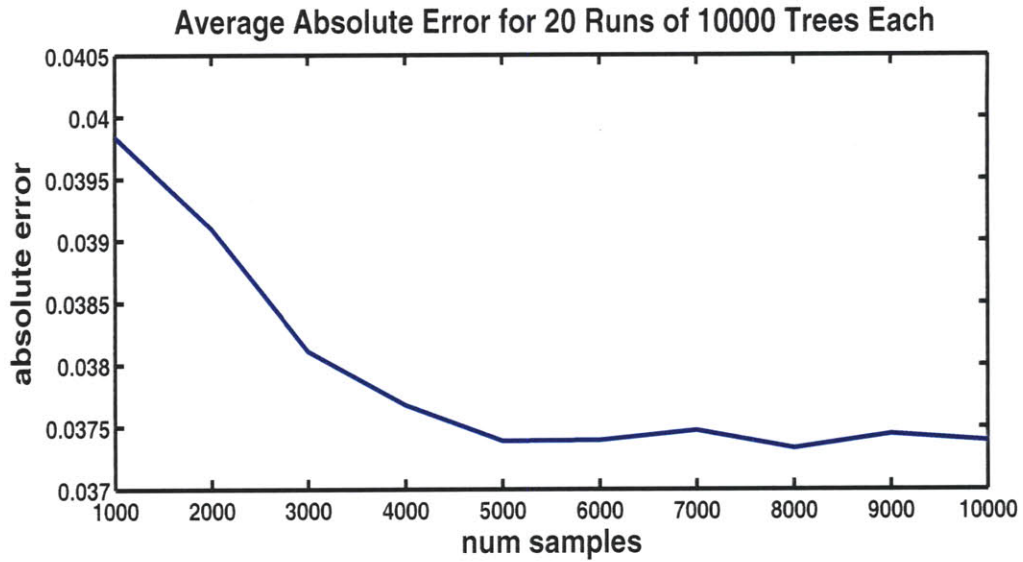


Figure 4-2: Results of run of Li *et al.*'s algorithm. As can be seen from the graphs above, both average absolute and squared error decreases as the number of samples increases.

that the algorithm was in fact converging to the desired distribution. Since time remained fixed, we could produce an exact discrete distribution on the possible set of trees. This distribution is shown in Figure 4-1 along with the marginal distribution of the 15 possible topologies. At every 1,000 sample increment, we computed both the average absolute error and average squared error between the distribution on topologies (determined by normalizing a histogram of all the topologies produced by the algorithm) for that number of samples per chain versus the actual distribution on topologies, and then averaged these errors across the 20 chains. The results are shown in Figure 4-2, and prove that the algorithm is in fact converging to the desired distribution.

For those wishing to implement the continuous-time version of the algorithm along with parameter re-estimation, the following points should be considered:

- Sampling time according to Equation (8) in [2] requires the use of Felsenstein's peeling algorithm [25]. Otherwise, to calculate the numerator of the density, we would need to sum 4^m terms (since the length of each sequence is m), an operation computationally infeasible for large data sets.
- Even with the use of the peeling algorithm, computing the the denominator of Equation (8) requires the use of numerical integration.
- Along the same lines, the procedure for re-estimating α and κ involves solving Equations (11) and (12) simultaneously, which once also requires the use of numerical methods (e.g., Newton-Raphson).

Chapter 5

Estimating Gene Phylogenies when Species Phylogeny is Known

Phylogeny can be reconstructed as described in Chapter 4 solely based on nucleotide sequence information. However, this model does not incorporate all of the information that is available for phylogenetic tree construction. In particular, the model developed by Li *et al.* [2] does not incorporate the known species tree information and the corresponding reconciliation (see Chapter 3), which provides a mapping from each gene in the gene tree to a species in the species tree and can aid in the search of the possible set of evolutionary trees. Rasmussen, in [26] showed that incorporating species tree information into a maximum likelihood framework for phylogenetic tree construction has potential for improving the accuracy of the respective algorithms. In this chapter, we formalize this idea by showing how to incorporate the additional species tree information into a model for a phylogenetic tree in a statistically rigorous way. We first formulate a new model which includes the species tree \mathcal{S}_{tree} and reconciliation \mathcal{R} as parameters and then propose a new Markov chain Monte Carlo algorithm for phylogenetic tree construction based on this model.

5.1 New Phylogenetic Tree Model

5.1.1 Model Parameters

We will denote the parameters in our model by $\theta = \{\mathcal{D}, \mathcal{T}, \mathcal{B}, \mathcal{S}_{tree}, \mathcal{R}\}$. The variables \mathcal{D} and \mathcal{T} are the same as that in the model in [2], which we have described in Chapter 4. To review, \mathcal{T} corresponds to the set of possible distinct labeled topologies of rooted bifurcating trees. In our case, each node in the tree represents a nucleotide sequence corresponding to a particular gene and we are interested in reconstructing the evolutionary history of this set of genes. For n external nodes, the number of topologies is given by $|\mathcal{T}| = \prod_{i=1}^{n-1} (2i - 1)$. Similarly, \mathcal{D} defines the set of nucleotide sequences for all of the nodes in the tree. Let D be the set of the four possible bases, i.e. $D = \{A, C, G, T\}$. We will denote each m -base nucleotide sequence by $\mathbf{d}_i \in D^m$ where $i \in \{-1, -2, \dots, -(n-1)\}$ for internal nodes (with -1 corresponding to the root) and $i \in \{1, 2, \dots, n\}$ for the external nodes. This means there are $|D|^{m(2n-1)}$ possible nucleotide sequences. Therefore the set of all sequences can conveniently be written as $\mathcal{D} = (\mathbf{d}_{-1}, \mathbf{d}_{-2}, \dots, \mathbf{d}_{-(n-1)}, \mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_n)$.

The first distinct feature of this model is that associated with each branch, with child node l , is a variable $b_l \in \mathbb{R}_+$ that governs how many substitutions, including back mutations, have occurred per site from the source of the branch (i.e., the parent) to the destination (i.e., the child). Recall that in the model in Chapter 4, branch lengths corresponded to evolutionary time. Collectively, we will refer to the set of branch lengths as $\mathcal{B} = (b_{-2}, b_{-3}, \dots, b_{-(n-1)}, b_1, b_2, \dots, b_n)$ ¹. Hence $B \in \mathbb{R}_+^{2n-2}$.

We also add two new variables into our model – the species tree \mathcal{S}_{tree} and the reconciliation mapping \mathcal{R} . Let the number of leaves in the species tree (representing modern species) be $S > 1$. Therefore the total number of nodes in the tree (modern and ancient species) is $2S - 1$. Each node S_j , where $j \in \{-1, -2, \dots, -(S-1)\}$ for internal nodes and $j \in \{1, 2, \dots, S\}$ for external nodes, corresponds to one of the $2S - 1$ possible modern or ancient species, such as human and mouse. Note in particular that the number of modern species S in the species tree need not equal

¹ b_{-1} is not included since there is no branch ending at the root node.

the number of leaves n in the gene tree. Each species tree is associated with a parameter g which corresponds to the overall mutation rate for the gene family represented by the tree. Associated with each branch in the species tree is a parameter $s_j \in \mathbb{R}_+$, where j is the index of the child node, that describes the species-specific mutation rate. This rate corresponds to a scaling of the gene-specific mutation rate specific to each species. Putting all these parameters together defines a species tree $\mathcal{S}_{tree} = (g, b_{-2}, b_{-3}, \dots, b_{-(S-1)}, b_1, b_2, \dots, b_S, S_{-1}, S_{-2}, \dots, S_{-(S-1)}, S_1, S_2, \dots, S_S)$ having dimension $4S - 2$. (i.e. one species parameter per node, one parameter per branch describing the species-specific mutation rate, and one parameter for the gene-specific rate).

The goal is to use the structure of the species tree to aid in the determination of the gene tree along with the already available sequence information. Thus, it is crucial to explicitly model the relationship between the two trees and their associated parameters. To this end, we introduce a new variable \mathcal{R} that encodes an explicit mapping from every node in the gene tree to some node in the species tree. If we denote the reconciliation of each gene in the gene tree by the index of the species node in the species tree, then we can represent the mapping \mathcal{R} as a vector $\mathcal{R} \in \{-1, -2, \dots, -(S-1), 1, 2, \dots, S\}^{2n-1}$. Since we associate each node with the branch coming from its parent, we will equivalently talk about the reconciliation as a mapping from branches in the gene tree to branches in the species tree.

Independence Assumptions Similar to the model in the Chapter 4, we pursue a Bayesian approach to modeling the statistical relationships among the variables described above. In order to arrive at a computationally tractable model, we impose a number of assumptions on the model. They are as follows:

- The topology \mathcal{T} is independent of the species tree \mathcal{S}_{tree} (i.e., $P(\mathcal{T}|\mathcal{S}_{tree}) = P(\mathcal{T})$).
- Given the branch lengths \mathcal{B} , the nucleotide sequences \mathcal{D} are conditionally independent of the species tree \mathcal{S}_{tree} and the reconciliation \mathcal{R} , since the branch lengths summarize all the information from the other two variables.

- In the gene tree \mathcal{T} , the mutation rates of the nucleotides in the sequence \mathbf{d}_i at node i are independent and identically distributed according to b_i^2 .

Given these assumptions the joint probability of all the parameters $\mathcal{B}, \mathcal{T}, \mathcal{D}, \mathcal{S}_{tree}$, and \mathcal{R} can be factored as follows:

$$\begin{aligned} P(\mathcal{B}, \mathcal{T}, \mathcal{S}_{tree}, \mathcal{R}, \mathcal{D}) &= P(\mathcal{D}|\mathcal{B}, \mathcal{T}, \mathcal{S}_{tree}, \mathcal{R})P(\mathcal{B}|\mathcal{T}, \mathcal{S}_{tree}, \mathcal{R})P(\mathcal{R}|\mathcal{T}, \mathcal{S}_{tree})P(\mathcal{T}|\mathcal{S}_{tree})P(\mathcal{S}_{tree}) \\ &= P(\mathcal{D}|\mathcal{B})P(\mathcal{B}|\mathcal{T}, \mathcal{S}_{tree}, \mathcal{R})P(\mathcal{R}|\mathcal{T}, \mathcal{S}_{tree})P(\mathcal{T})P(\mathcal{S}_{tree}). \end{aligned}$$

5.1.2 Model Specification

Now, we will describe how to determine each term of the joint probability above in turn, proceeding from right to left.

Prior probability of the species tree: $P(\mathcal{S}_{tree})$. The topology of the species tree and all associated parameters are fixed. Therefore, $P(\mathcal{S}_{tree})$ is known and a constant. The procedure by which the \mathcal{S}_{tree} is learned as described in [26].

Prior probability of topology: $P(\mathcal{T})$. We assume $P(\mathcal{T})$ is uniform over the possible $\prod_{i=1}^{n-1}(2i-1)$ topologies, i.e.,

$$P(\mathcal{T}) = \frac{1}{\prod_{i=1}^{n-1}(2i-1)}.$$

Conditional probability of reconciliation given species tree and topology: $P(\mathcal{R}|\mathcal{T}, \mathcal{S}_{tree})$. Similarly, once we know the species tree and topology, the reconciliation is determined analytically (see Section 3.2.2 for details) so $P(\mathcal{R}|\mathcal{T}, \mathcal{S}_{tree})$ is also known and a constant.

Conditional probability of branch lengths given topology, species tree, and reconciliation: $P(\mathcal{B}|\mathcal{T}, \mathcal{S}_{tree}, \mathcal{R})$. Once we have the topology, species tree, and reconciliation, we can determine the distribution on branch lengths as described below

²This has been shown to be supported by real biological data by Rasmussen in [26].

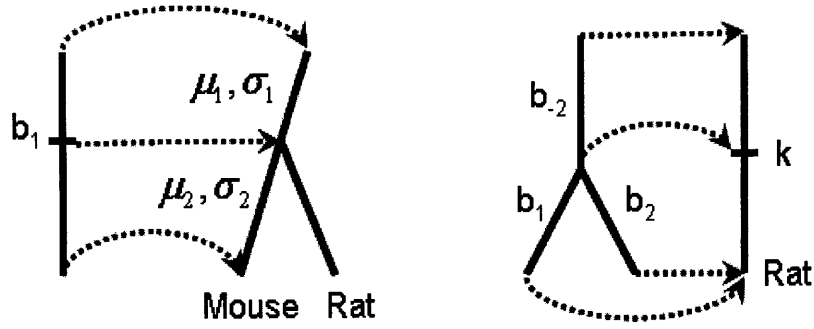


Figure 5-1: Two different complex reconciliations. On the left, one gene branch maps onto multiple species branches. On the right, multiple gene branches map onto a single species branch. Source: [26]

(we refer the reader to [26] for additional details): there are three possible types of reconciliations from a gene tree to a species tree.

The first type of reconciliation corresponds to the simple case in which a gene tree branch b_i reconciles to a unique species branch S_j . In this case, we model the distribution of the gene mutation rate conditioned on the species mutation rate by a normal distribution. In particular, $P(b_i|S_j) = \mathcal{N}(\mu_j, \sigma_j^2)$, where the parameters μ_j and σ_j^2 are known³.

The second type of reconciliation corresponds to the case where one gene branch maps onto multiple species branches and is illustrated on the left hand side of Figure 5-1. Fortunately, we can reduce this case to the simple type of reconciliation by merging the multiple species branches into one branch S_3 as shown in the figure. This merge recreates a one-to-one mapping of the gene branches to species branches. Since the sum of two independent Gaussian random variables is Gaussian, we have that: $P(b_3|S_3) = \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$. This idea easily extends to the case where we merge more than two species branches.

The last type of reconciliation arises in the case where multiple gene branches map onto one species branch S_j , as illustrated on the right hand side of Figure 5-1. In this case, we add points that split the species branch into multiple branches and once again produce a mapping where each branch maps to a single species branch

³In an alternative scheme, these parameters can be part of the model. We assume they are learned from data in the training phase as described in [26].

(which may be a part of one of the branches in the actual species tree). So, if μ_j and σ_j^2 are the parameters of the normal distribution associated with the parameter b_j of the actual species branch, we would like to have that $P(b_1 + b_{-2}|S_j) = \mathcal{N}(\mu_j, \sigma_j^2)$. We add in a point k to the species branch which ranges from 0 to 1 and represents the fraction of the species branch which is reconciled to a particular gene branch. Now we may set $P(b_1|S_j) = \mathcal{N}(k\mu_j, k\sigma_j^2)$ and $P(b_{-2}|S_j) = \mathcal{N}((1-k)\mu_j, (1-k)\sigma_j^2)$ thereby achieving our goal. To determine the likelihood of the three gene branches, we marginalize over k (we assume a uniform prior since we do not assume that any one value is more likely than any another). The equation takes the following form:

$$\begin{aligned} P(b_1, b_2, b_{-2}|S_j) &= \int_0^1 P(b_1, b_2, b_{-2}|S_j, k)P(k)dk \\ &= \int_0^1 P(b_1|S_j, k)P(b_2|S_j, k)P(b_{-2}|S_j, k)P(k)dk. \end{aligned}$$

The integral in the above equation can be computed using numerical integration. Gaussian quadrature was used in [26]. If we have more than two gene branches that reconcile to a single species branch, we need to include several breakpoints. Our likelihood then requires the evaluation of an integral for each breakpoint value which ranges uniformly from 0 to the previous breakpoint value. We do this, since once one breakpoint is chosen a gene further down the gene tree can only reconcile below that value because of the time ordering constraints imposed on the genes by the structure of the gene tree.

The complete distribution $P(\mathcal{B}|\mathcal{T}, \mathcal{S}_{tree}, \mathcal{R})$ is a product of the distributions across all branches, where the branch distributions are determined according to which of these three types of reconciliations is present. There is no closed form for this distribution, though we will illustrate what the distribution looks like for the case shown in Figure 5-2, where one of each of the three types of reconciliations is present. Let species branch S_j be parameterized by μ_j and σ_j specifying the distribution of the mutations rates associated with this branch. In the example in the figure, gene tree branch b_{-3} maps onto a single species branch S_4 . This is a simple type of reconciliation and hence the distribution $P(b_{-3}|S_4)$ is equal to $\mathcal{N}(\mu_4, \sigma_4^2)$.

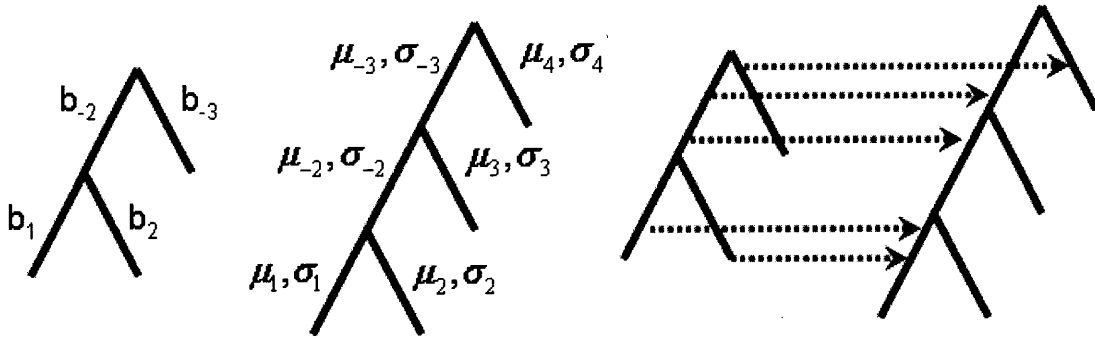


Figure 5-2: Reconciliation of a gene tree. The variables for the branch lengths and species tree parameters are shown on the left. The right-hand figure shows the mapping from gene tree to species tree branches. One of each of the three types of possible reconciliations is present.

Gene tree branch b_{-2} maps onto two species branches, S_{-3} and S_{-2} . This is the first type of complex reconciliation, in which a single gene tree branch maps onto multiple species tree branches. Hence, branch length b_{-2} is distributed according to the sum of two independent normal random variables, i.e. $P(b_{-2}|S_{-3}, S_{-2}) = \mathcal{N}(\mu_{-3} + \mu_{-2}, \sigma_{-3}^2 + \sigma_{-2}^2)$. Finally, branches b_1 and b_2 both map onto the same species branch S_1 . This is the second type of complex reconciliation in which two gene tree branches map onto a single species branch. Therefore, we need to add in a breakpoint k onto species branch S_1 and marginalize over it. This produces the following distribution on the two branch lengths: $P(b_1, b_2|S_1) = \int_0^1 \mathcal{N}(k\mu_1, k\sigma_1^2) \mathcal{N}((1-k)\mu_1, (1-k)\sigma_1^2) dk$. We can now get the complete distribution of all of the branch lengths in the gene tree by taking the product of these three distributions. Hence, $P(b_1, b_2, b_{-2}, b_{-3}|S_1, S_2, S_3, S_4, S_{-2}, S_{-3}, \mathcal{T}, \mathcal{R}) = \mathcal{N}(\mu_4, \sigma_4^2) \mathcal{N}(\mu_{-3} + \mu_{-2}, \sigma_{-3}^2 + \sigma_{-2}^2) * \int_0^1 \mathcal{N}(k\mu_1, k\sigma_1^2) \mathcal{N}((1-k)\mu_1, (1-k)\sigma_1^2) dk$.

Conditional probability of gene tree sequences given branch lengths: $P(\mathcal{D}|\mathcal{B})$.

Finally, only one term remains in the joint probability which we still need to compute: $P(\mathcal{D}|\mathcal{B})$. Once we know the per site mutation rate for each branch b_i , we can use this value to determine the probability that one sequence evolved to another sequence. We can use a similar idea to that in [2] and define $P_{\mathbf{d}_k \mathbf{d}_l}(b_l)$ as the probability that the nucleotide sequence at node k , denoted by \mathbf{d}_k , changes to the nucleotide sequence

at node l , denoted by \mathbf{d}_l , when the branch length between the two sequences is b_l . For the Jukes-Cantor model, the probability D_s that a single site mutation occurs, given a branch length b_l is given by the following formula [25]:

$$D_s = \frac{3}{4}(1 - e^{-\frac{4}{3}b_l}).$$

If we define the number of substitutions (transitions and transversions) between the two sequences as $tr_{\mathbf{d}_k \mathbf{d}_l}$, then, for a sequence of length m , we have that:

$$P_{\mathbf{d}_k \mathbf{d}_l}(b_l) = (D_s)^{tr_{\mathbf{d}_k \mathbf{d}_l}}(1 - D_s)^{m - tr_{\mathbf{d}_k \mathbf{d}_l}},$$

where we have used the fact that sites evolve independently to write the probability in terms of a product over the m sites.

As before, the probability of starting at the root can be computed as $\prod_{j \in D} \pi_j^{r_p}$, where r_p be the number of sites in the root node for which the nucleotide is p , where $p = A, C, G, \text{ or } T$. Since we have assumed that branches evolve independently, the entire conditional probability we are interested in computing is simply a product of starting at the root times the product of seeing all the transitions in the tree. We can write this quantity as follows:

$$P(\mathcal{D}|\mathcal{B}) = \left(\prod_{p \in D} \pi_j^{r_p} \right) \prod_{k, l \in E_{\mathcal{T}}} P_{\mathbf{d}_k \mathbf{d}_l}(b_l),$$

where $E_{\mathcal{T}}$ is the set of edges in \mathcal{T} .

5.2 Inference Algorithm

In this section, we describe Li *et al.*'s [2] algorithm for phylogenetic tree construction. We begin with a motivation for developing the algorithm and then proceed to the algorithm's description.

5.2.1 Motivation for Developing an MCMC-based Inference Algorithm

For a given problem, the observed variables are the nucleotide sequences for the external nodes, which we denote as $\mathcal{D}_{ext} = (\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_n)$ ⁴ as well as the species tree \mathcal{S}_{tree} . We denote the rest of the nucleotide sequences by $\mathcal{D}_{int} = (\mathbf{d}_{-1}, \mathbf{d}_{-2}, \mathbf{d}_{-3}, \dots, \mathbf{d}_{-(n-1)})$ so that $\mathcal{D} = \mathcal{D}_{int} \cup \mathcal{D}_{ext}$. The quantity that we are interested in computing is $E[\mathcal{B}, \mathcal{T}, \mathcal{R}, \mathcal{D}_{int} | \mathcal{D}_{ext}, \mathcal{S}_{tree}]$, which requires the computation of the following conditional probability distribution:

$$P(\mathcal{B}, \mathcal{T}, \mathcal{R}, \mathcal{D}_{int} | \mathcal{D}_{ext}, \mathcal{S}_{tree}) = \frac{P(\mathcal{B}, \mathcal{T}, \mathcal{S}_{tree}, \mathcal{R}, \mathcal{D}_{int}, \mathcal{D}_{ext})}{\sum_{\mathcal{T}} \sum_{\mathcal{D}} \sum_{\mathcal{B}} P(\mathcal{B}, \mathcal{T}, \mathcal{S}_{tree}, \mathcal{R}, \mathcal{D}_{int}, \mathcal{D}_{ext})}. \quad (5.1)$$

We can rewrite the joint probability in the numerator and denominator using the conditional distributions we have previously described. Unfortunately, since $P(\mathcal{B} | \mathcal{T}, \mathcal{S}_{tree}, \mathcal{R})$ does not have a closed form, we cannot compute the normalizing constant directly. Even if we employed a method that approximated this term, the computation of the normalizing constant in the denominator would be expensive even for small values of n . This problem therefore lends itself nicely to the application of MCMC methods. We can use MCMC to draw samples from the distribution in Equation (5.1) and then compute the expectation we are interested in by taking the average of the samples we generate. The Strong Law of Large Numbers guarantees that the sample mean will converge to the true mean as the number of samples increases. In the following section, we describe an MCMC procedure for obtaining the desired estimates.

5.2.2 Algorithm Description

We will now outline a new algorithm for phylogenetic tree construction aimed at the computation of our quantity of interest, namely $E[\mathcal{B}, \mathcal{T}, \mathcal{R}, \mathcal{D}_{int} | \mathcal{D}_{ext}, \mathcal{S}_{tree}]$. The algorithm will utilize nearest-neighbor interchange as the local rearrangement strat-

⁴Note: this is different than in the algorithm in Chapter 4 in which the root node sequence was also observed.

egy based on the model we have previously described, which incorporates species tree information and the reconciliation. To compute our quantity of interest, we define a Markov chain whose stationary distribution is the target distribution $\pi = P(\mathcal{B}, \mathcal{T}, \mathcal{R}, \mathcal{D}_{int} | \mathcal{D}_{ext}, \mathcal{S}_{tree})$. The pseudocode for the algorithm is shown in Algorithm 8 and a detailed description is provided below.

Algorithm 8: New Algorithm

Data: $\mathcal{D}_{ext}, \mathcal{S}_{tree}$
Result: N trees $\{\mathbf{X}_i; i = 1, 2, \dots, N\}$ whose limiting distribution is π
begin
 Initialize \mathbf{X}_0 using neighbor-joining and reconciliation-rerooting [36]
 for $i = 1$ *to* N **do**
 Sample a tree \mathbf{Y}_i from $q(\mathbf{y} | \mathbf{X}_{i-1})$ as follows:
 1. Perform nearest-neighbor interchange to produce a new tree \mathbf{Y}_i
 2. Compute the reconciliation of \mathbf{Y}_i
 3. Determine new branch lengths for \mathbf{Y}_i
 4. Determine new sequence for target node T in \mathbf{Y}_i
 Compute $\alpha(\mathbf{X}_i, \mathbf{Y}_i)$ according to Equation (5.3)
 Sample a uniform random variable \mathbf{U} on $[0, 1]$
 if $\mathbf{U} \leq \alpha(\mathbf{X}_i, \mathbf{Y}_i)$ **then**
 | $\mathbf{X}_{i+1} = \mathbf{Y}_i$
 else
 | $\mathbf{X}_{i+1} = \mathbf{X}_i$
 end

We generate an initial tree using neighbor-joining and reconciliation-rerooting [36] (since neighbor-joining does not always produce a rooted tree, we need to perform a rerooting procedure). Once we have an initial tree, we propose new trees using the nearest-neighbor interchange (NNI) method, described in Section 3.2.5. In our case, since the tree is rooted, NNI works as follows: the rearrangement occurs in the neighborhood of one of the internal nodes T (called the target). Note that NNI for rooted trees is, in essence, the local rearrangement strategy proposed by Li *et al.* [2] with the exception that the root can be chosen as the target node. This neighborhood includes the target’s parent P , sibling S , and children $C1$ and $C2$. The algorithm based on the local rearrangement strategy can be described in the following steps:

- **Step 1:** Randomly choose a node among the $n - 1$ internal nodes to be the

target node T . From among $C1$, $C2$, and S , randomly select one of the nodes to be the new sibling S' . The remaining two nodes become the new children - $C1'$ and $C2'$.

- **Step 2:** Given the topology selected in Step 1, compute the reconciliation for the new gene tree to the known species tree using the algorithm describe in Section 3.2.2.
- **Step 3:** Once we know the topology and the reconciliation, we can sample new branch lengths for the following 4 branches: the branch going from the parent to the target node, the two branches going from the target node to the new children, and the branch going from the parent node to the new sibling. For each of the branches, we determine which of the three types of reconciliations is present, and sample the new branch length from the appropriate distribution. The three types of reconciliations are listed below:

Simple Reconciliation: Gene tree branch b_i matches to a unique species branch S_j . In this case, the new branch length $b'_i \sim N(\mu_j, \sigma_j^2)$.

Complex Reconciliation 1: In this case, one gene branch maps onto multiple species branches. The new branch length b'_i is sampled from the sum of the independent normals corresponding to the multiple species branches being mapped to. For the case of b'_i mapping to two species branches S_1 and S_2 , $b'_i \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$.

Complex Reconciliation 2: In this case, multiple gene branches map onto one species branch S_j . We split the species branch to accommodate all of the gene branches and then sample each new gene branch b'_i from the partial species branch. When there are two gene tree branches b'_1 and b'_2 that map onto a single species branch S_1 , we sample $k \sim \mathcal{U}(0, 1)$ and then $b'_1 \sim N(k\mu_1, k\sigma_1^2)$ and $b'_2 \sim N((1 - k)\mu_1, (1 - k)\sigma_1^2)$.

- **Step 4:** Having picked the new branch lengths, the last step of the algorithm is determining a new nucleotide sequence for the new target node, which we will

denote by $\mathbf{d}_{T'}$. The distribution on this sequence is as follows:

$$h(\mathbf{d}_{T'}) = \frac{P_{\mathbf{d}_P \mathbf{d}_{T'}}(b_{T'}) P_{\mathbf{d}_{T'} \mathbf{d}_{C1'}}(b_{C1'}) P_{\mathbf{d}_{T'} \mathbf{d}_{C2'}}(b_{C2'})}{\sum_{\mathbf{d} \in D^m} P_{\mathbf{d}_P \mathbf{d}}(b_T) P_{\mathbf{d} \mathbf{d}_{C1'}}(b_{C1'}) P_{\mathbf{d} \mathbf{d}_{C2'}}(b_{C2'})} \text{ for } \mathbf{d}_{T'} \in D^m \quad (5.2)$$

However, since we have made an independent site evolution assumption in our model, choosing a new sequence for a target node can be done by choosing each site in the sequence independently.

5.2.3 Computing the Acceptance Probability

To complete our algorithm description, we need to describe how to compute the acceptance probability for the tree produced after the local rearrangement strategy. As we have done in Chapter 4, if we denote our initial tree by \mathbf{X}_i and the tree produced after the rearrangement procedure described above as \mathbf{Y}_i , then the Markov chain transition probability density for the local rearrangement strategy is given by $P(\mathbf{X}_i, \mathbf{Y}_i) = \alpha(\mathbf{X}_i, \mathbf{Y}_i)q(\mathbf{Y}_i|\mathbf{X}_i)$, with

$$\alpha(\mathbf{X}_i, \mathbf{Y}_i) = \min \left(\frac{\pi(\mathbf{Y}_i)q(\mathbf{X}_i|\mathbf{Y}_i)}{\pi(\mathbf{X}_i)q(\mathbf{Y}_i|\mathbf{X}_i)}, 1 \right), \quad (5.3)$$

$$\pi(\mathbf{X}_i) = P(\mathcal{B}, \mathcal{T}, \mathcal{R}, \mathcal{D}_{int} | \mathcal{D}_{ext}, \mathcal{S}_{tree}) = P(\mathcal{B}, \mathcal{T}, \mathcal{R}, \mathcal{D} \mathcal{S}_{tree}) P(\mathcal{D}_{ext}, \mathcal{S}_{tree}) \quad (5.4)$$

$$= P(\mathcal{B}, \mathcal{T}, \mathcal{R}, \mathcal{D} \mathcal{S}_{tree}) P(\mathcal{D}_{ext}) P(\mathcal{S}_{tree}) = P(\mathcal{B}, \mathcal{T}, \mathcal{S}_{tree}, \mathcal{R}, \mathcal{D}) \quad (5.5)$$

$$= P(\mathcal{D} | \mathcal{B}) P(\mathcal{B} | \mathcal{T}, \mathcal{S}_{tree}, \mathcal{R}) P(\mathcal{R} | \mathcal{T}, \mathcal{S}_{tree}) P(\mathcal{T}) P(\mathcal{S}_{tree}), \quad (5.6)$$

$$\pi(\mathbf{Y}_i) = P(\mathcal{D}' | \mathcal{B}') P(\mathcal{B}' | \mathcal{T}', \mathcal{S}_{tree}, \mathcal{R}') P(\mathcal{R}' | \mathcal{T}', \mathcal{S}_{tree}) P(\mathcal{T}') P(\mathcal{S}_{tree}), \quad (5.7)$$

$$q(\mathbf{Y}_i | \mathbf{X}_i) = P(\mathcal{T}' | \mathcal{T}) P(\mathcal{R}' | \mathcal{T}', \mathcal{S}_{tree}) P(\mathcal{B}' | \mathcal{T}', \mathcal{S}_{tree}, \mathcal{R}') P(\mathcal{D}'_{int} | \mathcal{D}_{int}, \mathcal{B}', \mathcal{T}', \mathcal{S}_{tree}, \mathcal{R}'), \quad (5.8)$$

$$q(\mathbf{X}_i | \mathbf{Y}_i) = P(\mathcal{T} | \mathcal{T}') P(\mathcal{R} | \mathcal{T}, \mathcal{S}_{tree}) P(\mathcal{B} | \mathcal{T}, \mathcal{S}_{tree}, \mathcal{R}) P(\mathcal{D}_{int} | \mathcal{D}'_{int}, \mathcal{B}, \mathcal{T}, \mathcal{S}_{tree}, \mathcal{R}), \quad (5.9)$$

where ' signifies a post local rearrangement parameter.

First, note that the bulk of the computation of the acceptance probability lies

in determining the quantities $\pi(\mathbf{X}_i)$ and $\pi(\mathbf{Y}_i)$, which correspond to the likelihoods of the trees represented by \mathbf{X}_i and \mathbf{Y}_i , respectively. We can calculate the likelihood for a tree, using the model we developed in Section 5.1.2, as follows: factor the joint probability density as shown in Equations (5.6) and (5.7) above, and then compute each term in turn as described in Section 5.1.2.

Now, we describe how to compute the quantities $q(\mathbf{Y}_i|\mathbf{X}_i)$ and $q(\mathbf{X}_i|\mathbf{Y}_i)$, the knowledge of which will allow us to determine the acceptance probability of the new tree. To determine $q(\mathbf{Y}_i|\mathbf{X}_i)$, we need to compute all of the terms in Equation (5.8). Proceeding from left to right, we first have the probability of the new topology \mathcal{T}' given the old topology \mathcal{T} . If the root node were chosen as the target in Step 1 of the algorithm, nearest-neighbor interchange (NNI) does not produce a new topology (since swapping the two child nodes of the root produces the same tree). Therefore, in this case, the term $P(\mathcal{T}'|\mathcal{T}) = 1$. Otherwise, if any of the other internal nodes were chosen as the target, there are three equally likely tree topologies post NNI, resulting in $P(\mathcal{T}'|\mathcal{T}) = \frac{1}{3}$. The second term $P(\mathcal{R}'|\mathcal{T}', \mathcal{S}_{tree})$, or the probability of the new reconciliation, given the species tree and the new topology, is a constant 1, since the new reconciliation is determined analytically. The third term $P(\mathcal{B}'|\mathcal{T}', \mathcal{S}_{tree}, \mathcal{R}')$, which is the probability of the new branch lengths, given the new topology, the species tree, and the new reconciliation, can be determined using the procedure described in Section 5.1.2. We use the new reconciliation to determine the joint distribution on branch lengths, and evaluate the likelihood of our proposed branch lengths given this joint distribution. The final term, $P(\mathcal{D}'_{int}|\mathcal{D}_{int}, \mathcal{B}', \mathcal{T}', \mathcal{S}_{tree}, \mathcal{R}')$, corresponds to the probability of the new target node sequence given the sequences of the old topology, the new branch lengths, the new topology, the species tree, and the new reconciliation. To compute this term, we use Equation (5.2) (i.e., count up the number of mutations that have occurred between the relevant sequences and then plug these values into the equation). The other conditional probability, $q(\mathbf{X}_i|\mathbf{Y}_i)$, can be computed in exactly the same way as described above with the primes switched (i.e., we instead condition on the new values of the parameters). Finally, we can put these two pieces together and compute the acceptance probability in Equation (5.3).

5.3 Implementation and Results

Although we have not implemented the algorithm described in Algorithm 8, an independent implementation of a modified version of this algorithm, which computes the maximum likelihood estimate based on maximizing the likelihood defined by the product of the two terms $P(\mathcal{D}|\mathcal{B})$ and $P(\mathcal{B}|T, \mathcal{S}_{tree}, \mathcal{R})$ in the model described in Section 5.1, has shown very promising results. In [26], Rasmussen performed tests for which he knew the correct tree topology and then measured the accuracy of several existing algorithms, including Mr. Bayes, in how well they were able to reconstruct the known topology. Our algorithm “had the highest accuracy relative to all of the standard phylogenetic tree construction algorithms described in [26]” (Matt Rasmussen, MIT Computational Biology Group, August 20, 2007). This initial result is promising and serves as further motivation for the use of the algorithm described in this chapter.

Chapter 6

Conclusion

Chapter 2 served as a stand-alone introduction to Markov chain Monte Carlo (MCMC). In Chapter 3, we provided an overview of the field of phylogenetics, which included a description of popular distance and model-based approaches to phylogenetic tree construction. Chapter 4 described the work of Li *et al.* [2], who developed a model and algorithm for reconstructing species phylogenies based on nucleotide sequences. Their work and the work of Rasmussen [26] was the basis for the model and algorithm we developed in Chapter 5, which incorporates known species phylogeny in order to reconstruct gene phylogenies.

The new algorithm we have developed shows promise for outperforming existing algorithms for phylogenetic tree construction. By incorporating species tree information into the model of a phylogenetic tree, we are able to bias the likelihood of proposed phylogenetic trees so that ones implying large numbers of gene duplications and losses are less likely. This fact may help improve the accuracy of phylogenetic tree construction algorithms and overcome the problems faced by other MCMC-based algorithms, including the one developed in [2].

A major contribution of this work, therefore, is the formulation of a new model which allows for the incorporation of the species phylogeny into phylogenetic tree construction process in a statistically rigorous way. Another contribution is the development of a phylogenetic tree construction algorithm based on this model. Finally, we have also provided a self-contained introduction to MCMC.

There are several possible extensions to consider for the new model we have developed. One involves re-estimating the parameters for the underlying sequence evolution model along with the gene-specific mutation rate during the burn-in period, based on the approach in [2]. This will be necessary as initial estimates are likely to be wrong. We would also like to extend our method to allow for other underlying sequence evolution models, as currently the algorithm is limited to the model developed by Jukes-Cantor. Another extension is implementing various convergence acceleration techniques for MCMC, such as Metropolis-coupled Markov chain Monte Carlo (MC³) [37] and simulated annealing [22], to expedite the convergence of the Markov chain underlying our algorithm. One can also consider different tree proposal strategies, perhaps a mixture approach which uses both local tree rearrangements together with larger, more pronounced rearrangements to aid in the exploration of the possible space of trees. Finally, incorporating a non-uniform distribution on topologies generated based on known species trees may further improve the accuracy of our algorithm.

Appendix A

Proofs

A.1 Ergodicity Theorem for Discrete-Time Finite State Space Markov Chains

In Chapter 2, we stated that an ergodic discrete-time finite state space Markov chain converges to a unique stationary distribution π . We will prove the following more general theorem allowing for a discrete state space whose proof implies the previously stated theorem:

Theorem A.1.1. *An ergodic discrete-time Markov chain converges to a unique stationary distribution π . The limiting probability $\pi_j = \lim_{n \rightarrow \infty} P(X_n = j) = \lim_{n \rightarrow \infty} p_{ij}^n$ where p_{ij} specifies the transition probability, i.e., $p_{ij} = P(X_{n+1} = j | X_n = i) \forall$ states $i, j \in S$.*

Proof: What we wish to prove is that if we have an ergodic finite-state discrete-time Markov chain, then

$$\pi_j = \lim_{n \rightarrow \infty} p_{ij}^n > 0 \quad \forall j. \tag{A.1}$$

We will prove that if the above inequality holds, then π is a stationary distribution and there exists no other stationary distribution. To see why this is true, first we will prove that when Equation (A.1) holds, there exists at least one stationary distribution $\hat{\pi}$ and then that any stationary distribution must equal the limiting distribution π .

First, note that the following inequalities hold:

$$\sum_{j=0}^m p_{ij}^n \leq \sum_{j=0}^{\infty} p_{ij}^n = 1 \quad \forall m,$$

since all $p_{ij}^n \geq 0$. Taking the limit as $n \rightarrow \infty$, we get:

$$\lim_{n \rightarrow \infty} \sum_{j=0}^m p_{ij}^n = \sum_{j=0}^m \lim_{n \rightarrow \infty} p_{ij}^n = \sum_{j=0}^m \pi_j \leq 1 \quad \forall m,$$

implying that $\sum_{j=0}^{\infty} \pi_j \leq 1$.

Rewriting the transition probability at time $n + 1$ as follows:

$$p_{ij}^{n+1} = \sum_{k=0}^{\infty} p_{ik}^n p_{kj} \geq \sum_{k=0}^m p_{ik}^n p_{kj} \quad \forall m,$$

and taking the limit as $n \rightarrow \infty$ as before, we get:

$$\pi_j \geq \sum_{k=0}^m \pi_k p_{kj} \quad \forall m \Rightarrow \pi_j \geq \sum_{k=0}^{\infty} \pi_k p_{kj} \quad \forall m.$$

We would like to show that the inequality on the RHS is an equality. Assuming that the inequality is strict for some j implies that the following is true:

$$\sum_{j=0}^{\infty} \pi_j > \sum_{j=0}^{\infty} \sum_{k=0}^{\infty} \pi_k p_{kj} = \sum_{k=0}^{\infty} \pi_k \sum_{j=0}^{\infty} p_{kj} = \sum_{k=0}^{\infty} \pi_k,$$

which is a contradiction. Therefore, it must be the case that:

$$\pi_j = \sum_{k=0}^{\infty} \pi_k p_{kj} \quad j = 0, 1, 2, \dots$$

If we define $\hat{\pi}_j = \frac{\pi_j}{\sum_{k=0}^{\infty} \pi_k}$, then $\hat{\pi}$ is a stationary distribution by definition. Hence, at least one stationary distribution exists.

Second, we wish to show that if $\hat{\pi}$ were any stationary distribution then it must

equal the limiting distribution π . If X_0 were the starting distribution, then:

$$\hat{\pi}_j = P(X_n = j) = \sum_{i=0}^{\infty} P(X_n = j | X_0 = i) P(X_0 = i) = \sum_{i=0}^{\infty} p_{ij}^n \hat{\pi}_i, \quad (\text{A.2})$$

which implies that $\hat{\pi}_j \geq \sum_{i=0}^m p_{ij}^n \hat{\pi}_i$. Taking the limit of $\hat{\pi}_j$ first as $n \rightarrow \infty$ and then as $m \rightarrow \infty$, we get:

$$\hat{\pi}_j \geq \lim_{m \rightarrow \infty} \lim_{n \rightarrow \infty} \sum_{i=0}^m p_{ij}^n \hat{\pi}_i \geq \sum_{i=0}^{\infty} \pi_j \hat{\pi}_i = \pi_j.$$

To show that $\hat{\pi}_j \leq \pi_j$, we can use the second equality in (A.2) and the fact that $\pi_{ij}^n \leq 1$ to rewrite $\hat{\pi}_j$ as follows:

$$\hat{\pi}_j \leq \sum_{i=0}^m p_{ij}^n \hat{\pi}_i + \sum_{i=m+1}^{\infty} \hat{\pi}_i \quad \forall m.$$

Take the limit as $n \rightarrow \infty$, we get

$$\hat{\pi}_j \leq \lim_{n \rightarrow \infty} \left(\sum_{i=0}^m p_{ij}^n \hat{\pi}_i + \sum_{i=m+1}^{\infty} \hat{\pi}_i \right) \leq \sum_{i=0}^m \pi_j \hat{\pi}_i + \sum_{i=m+1}^{\infty} \hat{\pi}_i \quad \forall m.$$

Since $\sum_{i=0}^{\infty} \hat{\pi}_i = 1$, when we take the limit as $m \rightarrow \infty$ of the above equation, we have

$$\hat{\pi}_j \leq \sum_{i=0}^{\infty} \pi_j \hat{\pi}_i = \pi_j.$$

□

A.2 Counting Distinct Labeled Rooted Bifurcating Trees

Here we prove that the total number of distinct labeled topologies for a rooted bifurcating tree with n external nodes is given by $|\mathcal{T}| = \prod_{i=1}^{n-1} (2i - 1)$.

Theorem A.2.1. *The total number of distinct labeled topologies for a rooted bifurcating tree with n external nodes is given by $|\mathcal{T}| = \prod_{i=1}^{n-1} (2i - 1)$.*

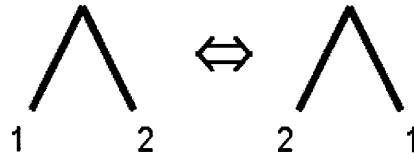


Figure A-1: Rooted bifurcating tree with two external nodes. If we swap nodes 1 and 2, we still have the same labeled topology.

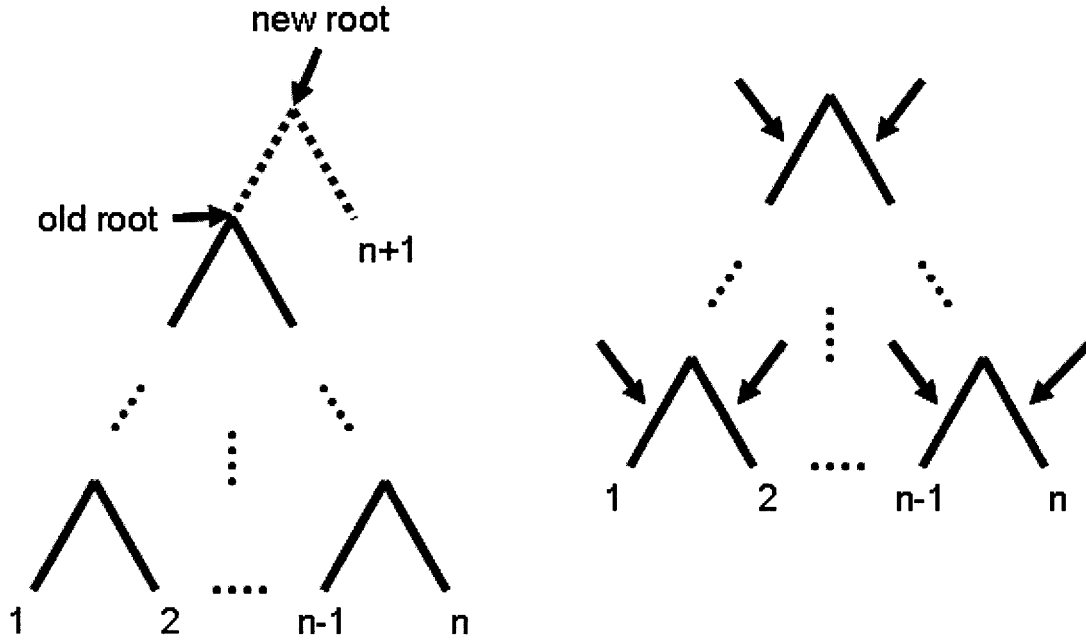


Figure A-2: Two ways of adding a new external node to a rooted bifurcating tree. We can either add it at the root, as shown on the left-hand side, or along any of the $2n - 2$ branches, as shown on the right-hand side.

Proof: We proceed by induction and first consider the base case when $n = 2$. This is the simplest bifurcating tree and if we label the two external nodes by 1 and 2, there is only one possible distinct labeled topology as shown in Figure A-1, since, if we swap 1 and 2, we still have the same labeled topology. This satisfies our formula since $\prod_{i=1}^{2-1} (2i - 1) = 1$.

Assuming the formula is true for n (i.e., our inductive hypothesis states that, for a rooted bifurcating tree with n external nodes, there are $\prod_{i=1}^{n-1}(2i - 1)$ distinct labeled topologies), we note that if we have n external nodes, there are two ways we can add in another external node. One possibility is to add it as a child of the root node as shown on the left-hand side of Figure A-2. In this case, by our inductive hypothesis, there are $\prod_{i=1}^{n-1}(2i - 1)$ possible ways to rearrange the other n external nodes. The other option is to add the last external node along any of the branches of the original tree, as shown on the right-hand side of figure. There are $2n - 2$ branches in a bifurcating rooted tree with n external nodes and hence $2n - 2$ possible places we can add the last external node. In each case, there are $\prod_{i=1}^{n-1}(2i - 1)$ ways to rearrange the other n external nodes. Enumerating the possibilities from the two options, we have a total of $\prod_{i=1}^{n-1}(2i - 1) + (2n - 2) \prod_{i=1}^{n-1}(2i - 1) = \prod_{i=1}^n(2i - 1)$ distinct labeled topologies. \square

Bibliography

- [1] M. Holder and P. Lewis, “Phylogeny estimation: traditional and Bayesian approaches,” *Nature Reviews Genetics*, vol. 4, pp. 275–284, 2003.
- [2] S. Li, J. Pearl, and H. Doss, “Phylogenetic tree construction using Markov chain Monte Carlo,” *Journal of the American Statistical Association*, vol. 95, pp. 493–508, 2000.
- [3] C. Robert and G. Casella, *Monte Carlo Statistical Methods*, Springer-Verlag, 2004.
- [4] W. R. Gilks, S. Richardson, and D. J. Spiegelhalter, *Markov chain Monte Carlo in Practice*, Chapman & Hall/CRC, 1996.
- [5] D. Higdon, “Auxiliary variable methods for Markov chain Monte Carlo with applications,” *Journal of the American Statistical Association*, vol. 93, pp. 585–595, 1998.
- [6] N. Lange, B. Carlin, and A. Gelfand, “Hierarchical Bayes models for the progression of HIV infection using longitudinal cd4+ t-cell numbers,” *Journal of the American Statistical Association*, vol. 87, pp. 615–626, 1992.
- [7] S. Biswas and S. Lin, “A Bayesian approach for incorporating variable rates of heterogeneity in linkage analysis,” *Journal of the American Statistical Association*, vol. 101, pp. 1341–1351, 2006.

- [8] G. Wei and A. Tanner, “A Monte Carlo implementation of the EM algorithm and the poor man’s data augmentation algorithms,” *Journal of the American Statistical Association*, vol. 85, pp. 699–704, 1990.
- [9] W. Hastings, “Monte Carlo sampling models using Markov chains and their applications,” *Biometrika*, vol. 57, pp. 97–109, 1970.
- [10] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller, “Equations of state calculations by fast computing machines,” *Journal of Chemical Physics*, vol. 21, pp. 1087–1092, 1953.
- [11] S. Geman and D. Geman, “Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 6, pp. 721–741, 1984.
- [12] G. Casella and E. George, “Explaining the Gibbs sampler,” *The American Statistician*, vol. 46, pp. 167–174, 1992.
- [13] M. Cowles and B. Carlin, “Markov chain Monte Carlo convergence diagnostics: a comparative review,” *Journal of the American Statistical Association*, vol. 91, pp. 883–904, 1996.
- [14] A. Gelman and D. Rubin, “Inference from iterative simulation using multiple sequences,” *Statistical Science*, vol. 7, pp. 457–472, 1992.
- [15] P. Mykland, L. Tierney, and B. Yu, “Regeneration in Markov chain samplers,” *Journal of the American Statistical Association*, vol. 90, pp. 233–241, 1995.
- [16] A. Gelfand and A.F.M. Smith, “Sampling-based approaches to calculating marginal densities,” *Journal of the American Statistical Association*, vol. 85, pp. 398–409, 1990.
- [17] C. J. Geyer, “Practical Markov chain Monte Carlo,” *Statistical Science*, vol. 7, pp. 473–511, 1992.

- [18] G.O. Roberts, A. Gelman, and W. Gilks, “Weak convergence and optimal scaling of random walk Metropolis algorithms,” *Annals of Applied Probability*, vol. 7, pp. 110–120, 1997.
- [19] R. Swendsen and J. Wang, “Nonuniversal critical dynamics in Monte Carlo simulations,” *Physical Review Letters*, vol. 58, pp. 86–88, 1987.
- [20] J. Besag and P. J. Green, “Spatial statistics and Bayesian computation,” *Journal of the Royal Statistical Society B*, vol. 55, pp. 25–37, 1993.
- [21] A. Gelfand and S. Sahu, “On Markov chain Monte Carlo acceleration,” *Journal of Computational and Graphical Statistics*, vol. 3, pp. 261–276, 1994.
- [22] C. J. Geyer and E. Thompson, “Annealing Markov chain Monte Carlo with applications to ancestral inference,” *Journal of the American Statistical Association*, vol. 90, pp. 909–920, 1995.
- [23] C. J. Geyer, “Markov chain Monte Carlo maximum likelihood,” in *Computing Science and Statistics: Proceedings of the 23rd Symposium on the Interface*, 1991, pp. 156–163.
- [24] S. Chib and E. Greenberg, “Understanding the Metropolis-Hastings algorithm,” *The American Statistician*, vol. 49, pp. 327–335, 1995.
- [25] J. Felsenstein, *Inferring Phylogenies*, Sinauer Associates, Inc., 2004.
- [26] M. Rasmussen, “Probabilistic framework for genome-wide phylogeny and orthology determination,” M.S. thesis, Massachusetts Institute of Technology, May 2006.
- [27] J. Eisen, “Phylogenomics: improving functional predictions for uncharacterized genes by evolutionary analysis,” *Genome Research*, vol. 8, pp. 163–167, 1998.
- [28] T. Jukes and C. Cantor, “Evolution of protein molecules,” *Mammalian Protein Metabolism*, vol. 3, pp. 21–132, 1969.

- [29] M. Kimura, “A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide sequences,” *Journal of Molecular Evolution*, vol. 16, pp. 111–120, 1980.
- [30] J. Felsenstein, “Evolutionary trees from DNA sequences: A maximum likelihood approach,” *Journal of Molecular Evolution*, vol. 17, pp. 368–376, 1981.
- [31] M. Hasegawa, T. Yano, and H. Kishino, “A new molecular clock of mitochondrial DNA and the evolution of hominoids,” *Proceedings of the Japanese Academy of Sciences B*, vol. 60, pp. 95–98, 1985.
- [32] R. D. M. Page, “Maps between trees and cladistic analysis of historical associations among genes, organisms, and areas,” *Systematic Biology*, vol. 43, pp. 58–77, 1994.
- [33] B. Mau, M. Newton, and B. Larget, “Bayesian phylogenetic inference via Markov chain Monte Carlo methods,” *Biometrics*, vol. 55, pp. 1–12, 1999.
- [34] B. Larget and D. Simon, “Markov chain Monte Carlo algorithms for the Bayesian analysis of phylogenetic trees,” *Journal of Molecular Biology and Evolution*, vol. 16, pp. 750–759, 1999.
- [35] Z. Yang and B. Rannala, “Bayesian phylogenetic inference using DNA sequences: a Markov chain Monte Carlo method,” *Journal of Molecular Biology and Evolution*, vol. 14, pp. 717–724, 1997.
- [36] C. Zmasek and S. Eddy, “RIO: analyzing proteomes by automated phylogenomics using resampled inference of orthologs,” *BMC Bioinformatics*, vol. 3, pp. 14, 2002.
- [37] G. Altekar, S. Dwarkadas, J. Huelsenbeck, and F. Ronquist, “Parallel Metropolis coupled Markov chain Monte Carlo for Bayesian phylogenetic inference,” *Bioinformatics*, vol. 20, pp. 407–415, 2004.