

# **Algorithms for Phylogenetic Reconstructions**

## **Lecture Notes**

Winter 2009/2010



# Preface

As the title implies, the focus of these notes is on ideas and algorithmic methods that are applied when evolutionary relationships are to be reconstructed from molecular sequences or other species-related data. Biological aspects like the molecular basis of evolution or methods of data acquisition and data preparation are not included.

A number of textbooks have been published covering most of the material discussed in these lecture notes, e.g. by Hillis *et al.* [1], Graur and Li [2], Gusfield [3], Page and Holmes [4], Felsenstein [5], and others. Nevertheless we believe this collection with a focus on algorithmic aspects of phylogenetic reconstruction methods is quite unique.

We thank Rainer Matthiesen for his notes taken during the Winter 2001/2002 lectures in Berlin, Rod Page for help with the practical course, and Heiko Schmidt for several discussions.

Berlin, February 2002

Martin Vingron, Jens Stoye, Hannes Luz

Meanwhile, this course has been given a number of times by various people: Jens Stoye (winter 2002/2003), Sebastian Böcker (winter 2003/2004), Marc Rehmsmeier (winter 2004/2005), Sven Rahmann and Constantin Bannert (winter 2005/2006). Each time the lecture notes were updated a little bit, and I have also made some additions for the winter 2006/2007 term. Nevertheless, some chapters are still not in a satisfactory state, sorry about this and all other errors contained.

Bielefeld, September 2006

Jens Stoye

---

After the last course (winter 2006/2007), Jens Stoye again updated the lecture notes. Meanwhile, Zsuzsanna Lipták gave a lecture with overlapping content and made some suggestions to improve these notes. During the preparation for the lecture in winter 2009/2010, I integrated her advices and contributed with own improvements and additions. Mainly, I revised the PP algorithm and the chapter on small parsimony. Thanks to Peter Husemann for his valuable comments and help.

Bielefeld, September 2009

Roland Wittler

# Contents

<b>I</b>	<b>Basics</b>	<b>1</b>
<b>1</b>	<b>Introduction: Creation, Evolution and Systematics</b>	<b>2</b>
<b>2</b>	<b>Graphs, Trees, and Phylogenetic Trees</b>	<b>6</b>
2.1	Graphs . . . . .	6
2.2	Trees . . . . .	7
2.3	Phylogenetic Trees . . . . .	8
2.3.1	What are Phylogenetic Trees? . . . . .	8
2.3.2	Tree Classification: Fully Resolved vs. Multifurcating . . . . .	11
2.3.3	Tree Representation and Tree Shape/Topology . . . . .	11
2.3.4	Weights in Phylogenetic Trees . . . . .	12
2.4	Counting Trees . . . . .	12
<b>3</b>	<b>Characters, States, and Perfect Phylogenies</b>	<b>14</b>
3.1	Characters and States . . . . .	14
3.2	Compatibility . . . . .	15
3.3	Perfect Phylogenies . . . . .	16
<b>II</b>	<b>Parsimony</b>	<b>21</b>
<b>4</b>	<b>The Small Parsimony Problem</b>	<b>22</b>
4.1	Introduction . . . . .	22
4.2	The Fitch Algorithm . . . . .	24
4.3	Hartigan's Algorithm . . . . .	26
4.4	The Sankoff Dynamic Programming Algorithm . . . . .	28
<b>5</b>	<b>Maximum Parsimony</b>	<b>30</b>
5.1	Introduction . . . . .	30
5.2	Exact Methods . . . . .	31
5.2.1	Brute-Force Method: Exhaustive Enumeration . . . . .	31
5.2.2	Branch and Bound Heuristic. . . . .	31

5.2.3	Greedy Sequential Addition . . . . .	33
5.3	Steiner Trees and Spanning Trees . . . . .	33
5.3.1	The DNA Grid Graph . . . . .	34
5.3.2	Steiner Trees . . . . .	34
5.3.3	Spanning Trees . . . . .	35
5.3.4	Spanning Trees and the Traveling Salesman Problem . . . . .	36
5.3.5	Spanning Trees as Approximations of Steiner Trees . . . . .	37
5.3.6	Application to Phylogeny . . . . .	37
5.4	Generalized Tree Alignment . . . . .	38
5.5	Long Branch Attraction . . . . .	38
<b>III</b>	<b>Distance-based Methods</b>	<b>39</b>
<b>6</b>	<b>Distance Based Trees</b>	<b>40</b>
6.1	Basic Definitions . . . . .	40
6.2	Ultrametric Trees . . . . .	43
6.2.1	Agglomerative Clustering . . . . .	43
6.3	Additive Trees . . . . .	46
6.3.1	Exact Reconstruction of Additive Trees . . . . .	47
6.3.2	Least Squares (Fitch-Margoliash) . . . . .	50
6.3.3	Minimum Evolution . . . . .	51
6.3.4	Neighbor Joining . . . . .	52
<b>7</b>	<b>Phylogenetic Networks</b>	<b>58</b>
7.1	Split Decomposition . . . . .	58
7.1.1	Introduction . . . . .	58
7.1.2	Basic Idea . . . . .	58
7.1.3	Definitions . . . . .	60
7.1.4	Computation of the $d$ -Splits . . . . .	61
7.2	NeighborNet . . . . .	63
<b>IV</b>	<b>Likelihood Methods</b>	<b>64</b>
<b>8</b>	<b>Modeling Sequence Evolution</b>	<b>65</b>
8.1	Basics on Probability . . . . .	65
8.1.1	Events and Probabilities . . . . .	65
8.1.2	Conditional Probability . . . . .	66
8.1.3	Bayes's Formula . . . . .	67
8.1.4	Independence . . . . .	67

8.1.5	Random Variables . . . . .	68
8.2	Markov Chains . . . . .	68
8.2.1	Time Discrete Markov Chains . . . . .	68
8.2.2	Time Continuous Markov Chains . . . . .	70
8.2.3	The Rate Matrix . . . . .	71
8.2.4	Definition of an Evolutionary Markov Process (EMP) . . . . .	72
8.3	Nucleotide Substitution Models . . . . .	74
8.3.1	The Jukes-Cantor Model (JC) . . . . .	74
8.3.2	Kimura 2-Parameter Model (K2P) . . . . .	76
8.4	Modeling Amino Acid Replacements . . . . .	76
8.4.1	Parameter Estimation . . . . .	77
8.4.2	Score Matrices . . . . .	77
8.4.3	Maximum Likelihood Estimation of Evolutionary Distances . . . . .	78
<b>9</b>	<b>Maximum Likelihood Trees</b>	<b>79</b>
9.1	Computing the Likelihood of a Given Tree . . . . .	79
9.2	Consistency . . . . .	82
<b>V</b>	<b>Advanced Topics</b>	<b>84</b>
<b>10</b>	<b>Assessing the Quality of Reconstructed Trees</b>	<b>85</b>
10.1	Bootstrapping . . . . .	85
10.1.1	The General Idea of Bootstrapping . . . . .	85
10.1.2	Bootstrapping in the Phylogenetic Context . . . . .	86
10.2	Likelihood Mapping . . . . .	86
<b>11</b>	<b>Consensus trees and consensus methods</b>	<b>88</b>
11.1	The Strict Consensus Tree . . . . .	88
11.2	Bayesian Methods . . . . .	89
<b>12</b>	<b>Parent trees and supertree methods</b>	<b>90</b>
12.1	Agreement Methods . . . . .	90
12.2	Optimization Methods . . . . .	90
12.3	One example: Quartet Puzzling . . . . .	90
<b>13</b>	<b>Fast Convergent Methods</b>	<b>92</b>
<b>14</b>	<b>The Method of Invariants</b>	<b>93</b>
14.1	Motivation: Development of a Consistent Method . . . . .	93
14.2	The Method of Invariants . . . . .	93



**Part I**  
**Basics**

# 1 Introduction: Creation, Evolution and Systematics

“Where do we come from” is one of the most frequently asked and discussed questions. Many theories and hypotheses suggest possible answers. One way to explain our origin is that we (and everything else) were created by an almighty god (Creationism, Intelligent Design). Often, creationism and evolution are presented as conflict, where people have to believe either in creation, or in evolution. However, we encourage another point of view: If we assume the existence of an almighty god, we can probably find answers to most questions. However, many of these answers depend on belief alone. In modern science, we try to answer our questions by experiments with reproducible results. This does not mean that we deny the existence of god, we just leave him/her/it out of consideration, because we can not prove or disprove his/her/its existence.

For centuries, biologists have tried to detect and classify the diversity in the biological world, this effort is known as *systematics*. If we assume that the biological diversity we see today is due to *evolution*, every species has a *phylogeny*, a history of its own evolutionary development. These two concepts gave rise to the science of *Phylogenetic Systematics*, where organisms are classified into groups by their phylogeny. Traditionally, this is based on morphological characters, while today molecular systematics prevail. Here is a short historical tour.

- Bible (*Genesis* 1, 26-27):

Then God said, “Let us make man in our image, in our likeness, and let them rule over the fish of the sea and the birds of the air, over the livestock, over all the earth, and over all the creatures that move along the ground.” So God created man in his own image, in the image of God he created him; male and female he created them.

- Carl Linné (1707-1778): Linné revolutionized the way in which species were classified and named. He proposed to group them by shared similarities into higher taxa, being: genera, orders, classes and kingdoms. He also invented the ‘binomial system’ for naming species, where a name is composed of genus and species in latin, as in *Homo sapiens*. The system rapidly became the standard system for naming species. In his early years as a researcher in botany, Linné believed in invariant species. Later on, he

---

admitted that a certain variation was possible. His most important works are *Systema naturae* (1735) and *Genera plantarum* (1737).

- Chevalier de Lamarck (1744-1829): Lamarck started as taxonomist in botany, and he was in general well educated. He applied Linné's ideas also to animals: *Philosophie zoologique* (1809). Lamarck was one of the first to believe in some kind of evolution. However, his way of explaining the observed changes was not accurate: In response to environmental changes, organisms would over- or underuse certain parts of their bodies. The heavily used parts would improve, the underused parts wither. These changes would be inherited by the offspring.
- Georges Cuvier (1769-1832): Cuvier's speciality was the comparison of organisms, characterizing their differences and similarities. He introduced the idea of the Phylum, but still believed in creation. He also studied fossils and observed changes in comparison with contemporary organisms. He explained them by a series of catastrophes. Each would have wiped out all life on earth, which would be newly created afterwards. Cuvier was a colleague of Lamarck, but he had no respect for Lamarck's theory of 'inheritance of acquired characteristics'.
- Charles Darwin (1809-1882): Here is an excerpt from the famous book *On the origin of species by means of natural selection* by Charles Darwin, 1859:

Whatever the cause may be of each slight difference in the offspring from their parents—and a cause for each must exist—it is the steady accumulation, through natural selection, of such differences, when beneficial to the individual, that gives rise to all the more important modifications of structure, by which the innumerable beings on the face of this earth are enabled to struggle with each other, and the best adapted to survive.

Darwin was not the first who believed in some sort of evolution. By the time he published his famous book, many biologists did not believe in the notion of fixed species anymore. However, Darwin was the one who was able to explain *how* this evolution could have occurred. His concept of evolution by natural selection may be expressed as a very simple set of statements:

1. The individual organisms in a population vary.
2. They overproduce (if the available resources allow).
3. Natural selection favors the reproduction of those individuals that are best adapted to the environment.
4. Some of the variations are inherited to the offspring.
5. Therefore organisms evolve.

- Alfred Russel Wallace (1823-1913) found during his extensive field work, besides others, the “Wallace line”, which separates Indonesia into two areas. In one part, he found species of Australian origin whereas in the other part, most species had Asian background. He and Darwin had similar ideas about *natural selection* and had continuous contact and discussions about their theories. Darwin decided to publish quickly and Wallace became a co-discoverer in the shadow of the more famous and acknowledged Darwin. Nevertheless, he settled for the resulting benefit. In contrast to Darwin, he explained the emergence of striking coloration of animals as a warning sign for adversaries and he devised the *Wallace effect*: Natural selection could inhibit hybridization and thus encourage speciation. Wallace also was an advocacy of spiritualism. In opposition to his colleagues, he claimed that something in the “unseen universe of Spirit” had interceded a least three times during evolution [6, p. 477].
- Ernst Haeckel (1834-1919) did a lot of field work and is known for his “genealogical tree” (1874, see Figure 1.1). He proposed the *biological law* “ontogeny recapitulates phylogeny”, which claims that the development of an individual reflects its entire evolutionary history.
- Emil Hans Willi Hennig (1913-1976) was specialized in *dipterans* (ordinary flies and mosquitoes). He noticed that morphological similarity of species does not imply close relationship. Hennig not only called for a phylogeny based systematic, but also stated corresponding problems, developed first, formal methods, and introduced an essential terminology. This was the basis for the *parsimony principle* and modern cladistics.
- Emil Zuckerkandl (\*1922) and Linus Pauling (1901-1994) were among the first to use biomolecular data for phylogenetic considerations. In 1962, they found that the number of amino acid differences in hemoglobin directly corresponds to time of divergence. They further abstracted from this finding and formulated the *molecular clock hypothesis* in 1965.

PEDIGREE OF MAN.

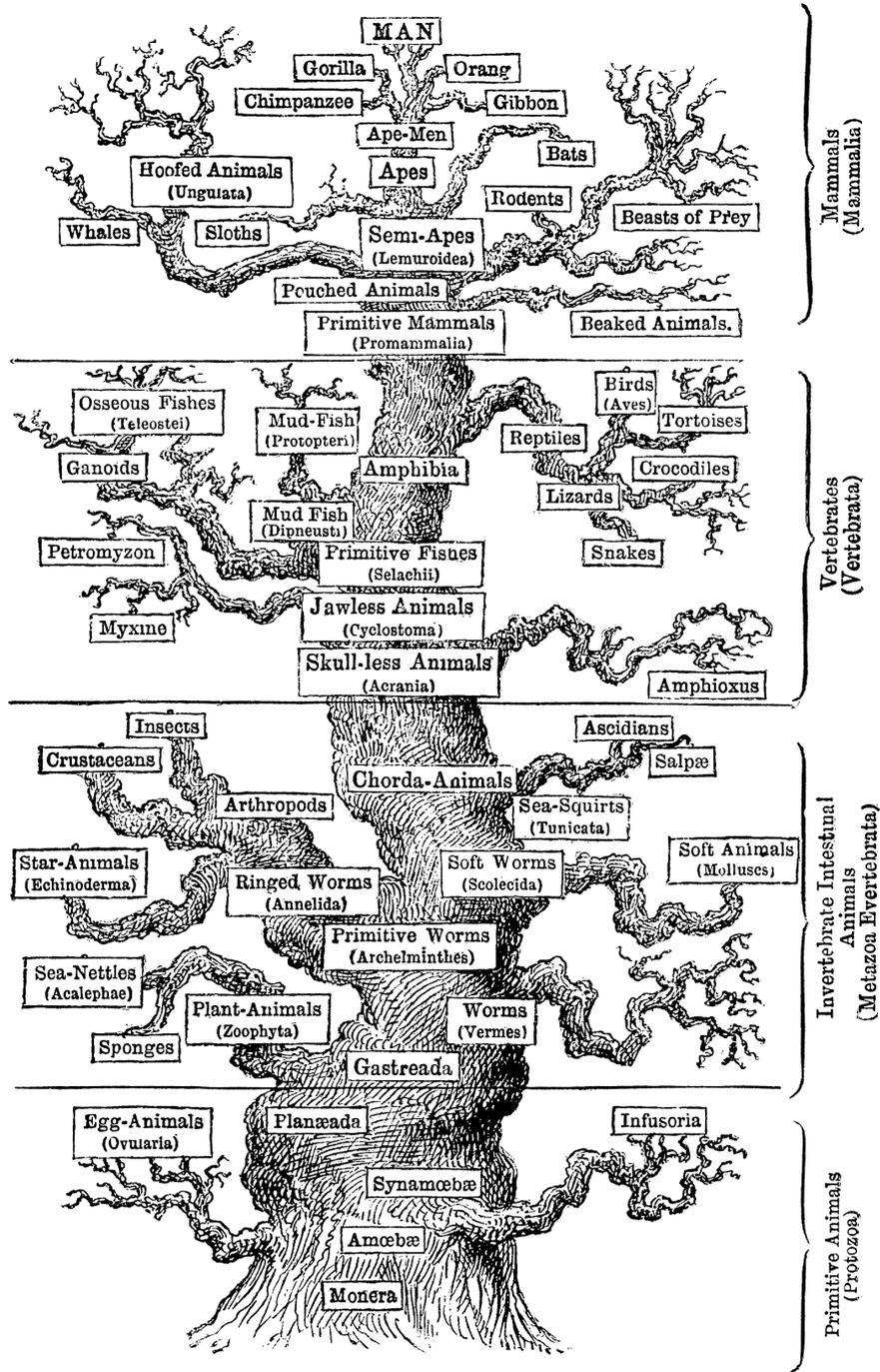


Figure 1.1: The genealogical tree by Ernst Haeckel, 1874.

## 2 Graphs, Trees, and Phylogenetic Trees

### 2.1 Graphs

In this course we talk a lot about trees. A tree is a special kind of graph. So we start with the definitions of a few graph theoretical terms before we see how trees are defined.

#### Undirected Graphs.

- An *undirected graph* is a pair  $G = (V, E)$  consisting of a set  $V$  of *vertices* (or *nodes*) and a set  $E \subseteq \binom{V}{2}$  of *edges* (or *branches*) that connect nodes. The number of nodes  $|V|$  is also called the *size* of  $G$ .
- The set  $\binom{V}{2}$  referred to above is the set of all 2-element subsets of  $V$ , i.e., the set of all  $\{v_1, v_2\}$  with  $v_1 \neq v_2$ . Sets are used to model undirected edges because there is no order among the vertices for an undirected edge.
- If  $e = \{v_1, v_2\} \in E$  is an edge connecting vertices  $v_1$  and  $v_2$ ,  $e$  is said to be *incident* to  $v_1$  and  $v_2$ . The two vertices  $v_1$  and  $v_2$  are *adjacent*.
- The *degree* of a vertex  $v$  is the number of edges incident to  $v$ .
- A *path* is a sequence of nodes  $v_1, v_2, \dots, v_n$  where  $v_i$  and  $v_{i+1}$  are connected by an edge for all  $i = 1, \dots, n - 1$ . The *length* of such a path is the number of edges along the path,  $n - 1$ . In a *simple path* all vertices except possibly the first and last are distinct. Two vertices  $v_i$  and  $v_j$  are *connected* if there exists a path in  $G$  that starts with  $v_i$  and ends with  $v_j$ . A graph  $G = (V, E)$  is *connected* if every two vertices  $v_i, v_j \in V$  are connected.
- A *cycle* is a path in which the first and last vertex are the same. A *simple cycle* uses each edge  $\{u, v\}$  at most once. A graph without simple cycles is called *acyclic*.

**Directed Graphs.**

- A *directed graph* (digraph) is a pair  $G = (V, E)$ , as for an undirected graph, but where now  $E \subseteq V \times V$ .
- An edge  $e = (v_1, v_2)$  is interpreted to point from  $v_1$  to  $v_2$ , symbolically also written as  $v_1 \xrightarrow{e} v_2$ . We usually assume that  $v_1 \neq v_2$  for all edges, although the definition does allow  $v_1 = v_2$ ; such an edge is called a *loop*. Graphs without loops are *loopless*.
- In a directed graph one distinguishes the *in-degree* and the *out-degree* of a vertex.

**Additional information in graphs.**

- If edges are annotated by numbers, one speaks of a *weighted graph*. Edge weights often represent lengths. The *length of a path* in a weighted graph is the sum of the edge lengths along the path.
- If edges are annotated by some label (e.g., letters or strings), one speaks of a *labeled graph*. One example is the definition of *suffix trees*.

Graphs are useful data structures in modeling real-world problems. Applications in bioinformatics include:

- modeling metabolic, regulatory, or protein interaction networks,
- physical mapping and sequence assembly (interval graphs),
- string comparison and pattern matching (finite automata, suffix trees),
- modeling sequence space,
- phylogenetic trees,
- modeling tree space.

## 2.2 Trees

A *tree* is a connected acyclic undirected graph. A *leaf* (*terminal node*) is a node of degree one. All other nodes are *internal* and have a degree of at least two. The *length of a tree* is the sum of all its edge lengths.

Let  $G = (V, E)$  be an undirected graph. The following statements are equivalent.

1.  $G$  is a tree.
2. Any two vertices of  $G$  are connected by a unique simple path.

3.  $G$  is minimally connected, i.e., if any edge is removed from  $E$ , the resulting graph is disconnected.
4.  $G$  is connected and  $|E| = |V| - 1$ .
5.  $G$  is acyclic and  $|E| = |V| - 1$ .
6.  $G$  is maximally acyclic, i.e., if any edge is added to  $E$ , the resulting graph contains a cycle.

One distinguishes between rooted and unrooted trees.

- An *unrooted tree* is a tree as defined above. An unrooted tree with degree three for all internal nodes is called a *binary tree*.
- A *rooted tree* is a tree in which one of the vertices is distinguished from the others and called the *root*. Rooting a tree induces a hierarchical relationships of the nodes and creates a directed graph, since rooting implies a direction for each edge (by definition always pointing away from the root). The terms *parent*, *child*, *sibling*, *ancestor*, *descendant* are then defined in the obvious way. Rooting a tree also changes the notion of the degree of a node: The *degree of a node in a rooted tree* refers to the *out-degree* of that node according to the above described directed graph. Then, a leaf is defined as a node of out-degree zero. A rooted tree with out-degree two for all internal nodes is called a *binary tree*. Each edge divides (splits) a tree into two connected components. Given a node  $v$  other than the root in a rooted tree, the *subtree rooted at  $v$*  is the remaining tree after deleting the edge that ends at  $v$  and the component containing the root. (The subtree rooted at the root is the complete, original tree.) The *depth* of node  $v$  in a rooted tree is the length of the (unique) simple path from the root to  $v$ . The *depth of a tree  $T$*  is the maximum depth of all of  $T$ 's nodes. The *width of a tree  $T$*  is the maximal number of nodes in  $T$  with the same depth.

## 2.3 Phylogenetic Trees

### 2.3.1 What are Phylogenetic Trees?

Phylogenetic (also: evolutionary) trees display the evolutionary relationships among a set of objects. Usually, those objects are species, but other entities are also possible. Here, we focus on species as objects, an example is shown in Figure 2.1 (left).

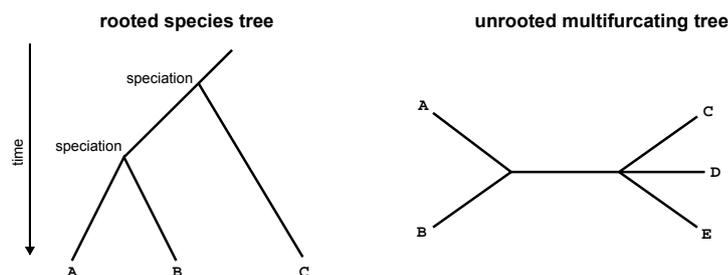
The  $n$  contemporary species are represented by the leaves of the tree. Internal nodes are branching (out-degree two or more in a rooted tree), they represent the last common ancestor before a speciation event took place. The species at the inner nodes are usually extinct, and

then the amount of data available from them is quite small. Therefore, the tree is mostly based on the data of contemporary species. It models their evolution, clearly showing how they are related via common ancestors.

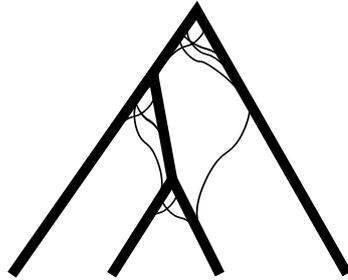
**Speciation.** The interpretation of phylogenetic trees requires some understanding of speciation, the origin of a new species. A speciation event is always linked to a population of organisms, not to an individual. Within this population, a group of individuals emerges that is able to live in a new way, at the same time acquiring a barrier to genetic exchange with the remaining population from which it arose (see, e.g. [7]). Usually, this is due to environmental changes that lead to a spatial separation, often by chance events. That is why speciation can not really be considered a punctual event, more realistic is a picture like in Figure 2.2.

After the separation of the two populations, both will diverge from each other during the course of time. Since both species evolve, the last common ancestor of the two will usually be extinct today, in the sense that the genetic pool we observe in either contemporary species is not the same as the one we would have observed in their last common ancestor.

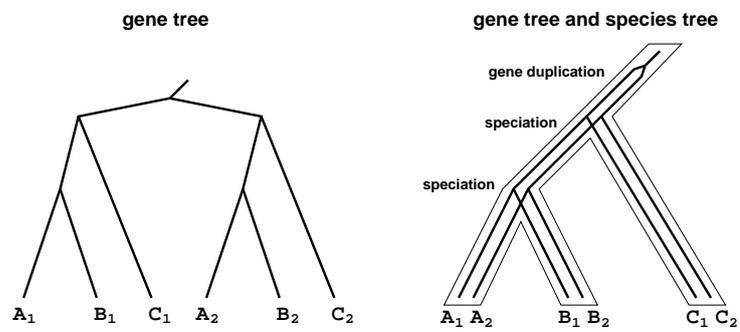
**Gene trees and species trees.** In modern molecular phylogeny, often the species at the leaves of a phylogenetic tree are represented by genes (or other stretches of genomic DNA). Such a tree is then often called a *gene tree*. In order to be comparable, all genes that are compared in the same tree should have originated from the same ancestral piece of DNA. Such genes are called *homologous genes* or *homologs*. It might also be that more than one homologous gene from the same species is included in the same tree. This can happen if the gene is duplicated and both copies are still alive. The full picture is then a mixture of speciations and gene duplications, see Figure 2.3. It is also possible that an internal node represents a *cluster* of species (the ones at the leaves in the subtree). Such a cluster is sometimes also called an *operational taxonomic unit* (OTU).



**Figure 2.1:** Left: Rooted, fully resolved species tree. Right: Unrooted tree with five taxa. The inner node branching to C, D and E is a polytomy.



**Figure 2.2:** A microscopic view at evolution.



**Figure 2.3:** Left: Gene tree, the same letter refers to the same species. Right: Species tree with a gene duplication.

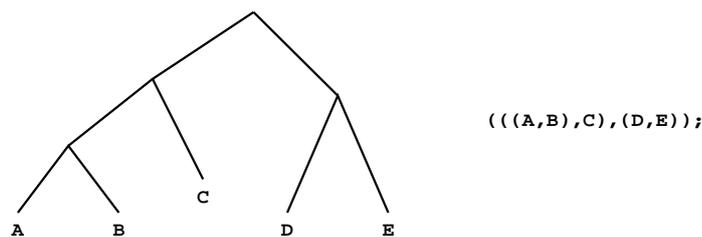
**Orthology and paralogy.** Two important concepts that are best explained with the help of gene trees and species trees are orthology and paralogy. Two genes are *orthologs* (or *orthologous genes*) if they descend from the same ancestor gene and their last common ancestor is a speciation event. They are *paralogs* (*paralogous genes*) if they descend from the same ancestor gene and their last common ancestor is a duplication event. For example, genes  $A_1$ ,  $B_1$  and  $C_1$  in Figure 2.3 are orthologs. In contrast, genes  $A_1$  and  $B_2$ , for example, are paralogs.

### 2.3.2 Tree Classification: Fully Resolved vs. Multifurcating

The difference between rooted and unrooted trees was already explained in Section 2.2. However, aside from being rooted or unrooted, Trees can also be *fully resolved* or not. For example, in a fully resolved rooted tree, each *inner* node has an in-degree of one, and an out-degree of two. The only exception is the root, which has an in-degree of zero. Trees that are not fully resolved (a.k.a. multifurcating trees) contain inner nodes with a degree of four or more. Those nodes are called unresolved nodes or *polytomies*, see Figure 2.1 (right). A polytomy can be due to simultaneous speciation into multiple lineages, or to the lack of knowledge about the exact speciation order.

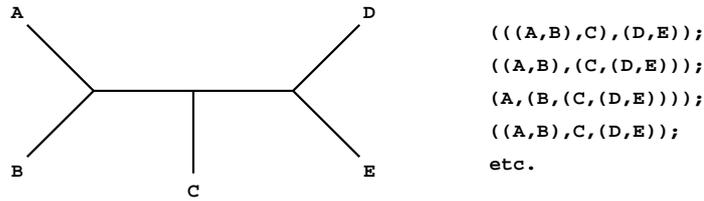
### 2.3.3 Tree Representation and Tree Shape/Topology

A simple and often used notation to represent a tree is the PHYLIP or NEWICK<sup>1</sup> format, where it is represented by a parenthesis structure, see Figure 2.4. This format implies a rooted tree where the root is represented by the outmost pair of parentheses. If unrooted trees are considered, different parenthesis structures can represent the same tree topology, see Figure 2.5.



**Figure 2.4:** A rooted tree with five leaves and the corresponding parenthesis structure.

<sup>1</sup><http://evolution.genetics.washington.edu/phylip/newicktree.html>



**Figure 2.5:** An unrooted tree and several parenthesis structures representing this tree.

It is important to notice that the branching order of edges at internal nodes is always arbitrary. Hence the trees  $(A, (B, C))$ ; and  $((B, C), A)$ ; are the same! (“Trees are like mobiles.”) Sometimes, as in nature, rooted trees are drawn with their root at the bottom. More often, though, rooted trees are drawn with their root at the top or from left to right. Unrooted trees are often drawn with their leaves pointing away from the center of the picture.

### 2.3.4 Weights in Phylogenetic Trees

Often in phylogeny one has *weighted* trees, where the edge lengths represent a distance. One example is the evolutionary time that separates two nodes in the tree. Here, trees are usually rooted, the time flowing from the root to the leaves. Another example for an edge weight is the number of morphological or molecular differences between two nodes. Generally, both rooted and unrooted weighted trees are considered in phylogeny. A special class of rooted weighted trees are those where each leaf has the same distance to the root, called *dendrograms*. Note that the same tree (same topology) can represent very different weighted trees!

Unlike general graphs, trees can always be drawn in the plane. However, sometimes it is difficult or impractical to draw the edge lengths of a tree proportional to their weight. In such cases one should write the edge lengths as numbers. However, one has to make sure that the reader does not confuse these numbers with other annotations like bootstrap values. (Bootstrap values are discussed in Chapter 10.)

## 2.4 Counting Trees

Finally we want to count how many different tree topologies exist for a given set of leaves. One has to distinguish between rooted and unrooted trees. Here we only consider unrooted trees.

**Observation.** Each unrooted binary tree with  $n$  leaves has exactly  $n - 2$  internal nodes and  $2n - 3$  edges.

Based on these formulas, we can prove the following relation.

**Lemma.** Given  $n$  objects, there are  $U_n = \prod_{i=3}^n (2i - 5)$  labeled unrooted binary trees with these objects at the leaves.

**Proof.** For  $n = 1, 2$  recall that by definition the product of zero elements is  $U_1 = U_2 = 1$ , and this is the number of labeled, unrooted trees with  $n = 1$ , respectively  $n = 2$  leaves. For  $n \geq 3$  the proof is by induction:

Let  $n = 3$ . There is exactly  $U_3 = \prod_{i=3}^3 (2i - 5) = 1$  labeled unrooted tree with  $n = 3$  leaves.

Let  $n > 3$ . There are exactly  $n$  leaves and thus  $2n - 3$  edges. A new edge (creating the  $(n + 1)$ -th leaf) can be inserted in any of the  $2n - 3$  existing edges. Hence, by induction hypothesis,

$$\begin{aligned} U_{n+1} &= U_n \cdot (2n - 3) \\ &= \prod_{i=3}^n (2i - 5) \cdot (2(n + 1) - 5) \\ &= \prod_{i=3}^{n+1} (2i - 5). \end{aligned}$$

□

Hence the number of unrooted labeled trees grows super exponentially with the number of leaves  $n$ . This is what makes life difficult in phylogenetic reconstruction. The number of unresolved (possibly multifurcating) trees is even larger.

# 3 Characters, States, and Perfect Phylogenies

## 3.1 Characters and States

Given a group of species and no information about their evolution, how can we find out the evolutionary relationships among them? We need to find certain properties of these species, where the following must hold:

- We can decide if a species has this property or not.
- We can measure the quality or quantity of the property (e.g., size, number, color).

These properties are called *characters*. The actual quality or quantity of a character is called its *state*. Formally, we have the following:

**Definition.** A *character* is a pair  $C = (\lambda, S)$  consisting of a property name  $\lambda$  and an arbitrary set  $S$ , where the elements of  $S$  are called *character states*.

### Examples:

- The existence of a nervous system is a binary character. Character states are elements of  $\{False, True\}$  or  $\{0, 1\}$ .
- The number of extremities (arms, legs,...) is a numerical character. Character states are elements of  $\mathbb{N}$ .

- Here is an alignment of DNA sequences:

```
Seq1:  A  C  C  G  G  T  A
Seq2:  A  G  C  G  T  T  A
Seq3:  A  C  T  G  G  T  C
Seq4:  T  G  C  G  G  A  C
```

A nucleotide in a position of the alignment is a character. The character states are elements of  $\{A,C,G,T\}$ . For the above alignment, there are seven characters (columns of the alignment) and four character states.

The definition given for characters and states is not restricted to species. Any *object* can be defined by its characters. It may be considered as a vector of characters.

**Example:** Bicycle, motorcycle, tricycle and car are objects. The number of wheels and the existence of an engine are characters of these objects. The following table holds the character states:

	# wheels	existence of engine
bicycle	2	0
motorcycle	2	1
car	4	1
tricycle	3	0

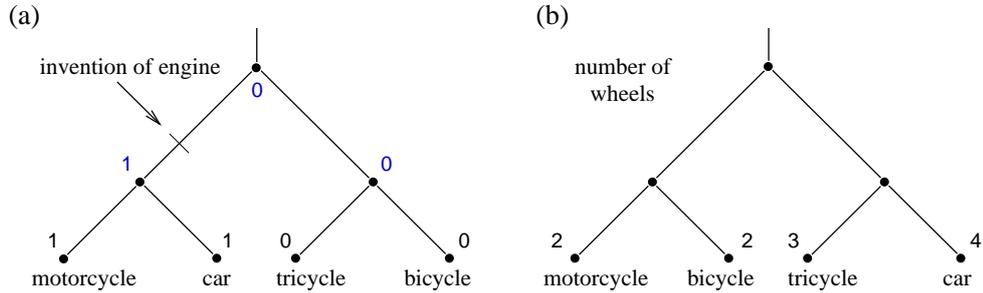
## 3.2 Compatibility

The terms and definitions of characters and states originate from paleontology. A main goal of paleontologists is to find correct phylogenetic trees for the species under consideration. It is generally agreed that the invention of a ‘new’ character state is a rare evolutionary event. Therefore, trees where the same state is invented multiple times independently are considered less likely in comparison with trees where each state is only invented once.

**Definition:** A character is *compatible* with a tree if all nodes of the tree can be labeled such that each character state induces one connected subtree.

Note: This implies that a character with  $k$  states that actually occur in our data will be compatible with a tree if we observe exactly  $k - 1$  changes of this state in that tree.

**Example:** Given a phylogenetic tree on a set of objects and a binary character  $c = \{0, 1\}$ : If the tree can be divided into two subtrees, where all nodes on one side have state 0, and 1 on the other, we count only one change of state. This can be seen in Figure 3.1 (a). The character ‘existence of engine’ is compatible with the tree, as the motor is invented only once (if the inner nodes are labeled correctly). The same character is not compatible with the tree in Figure 3.1 (b), where the engine is invented twice (if we assume that the common ancestor had no engine). The character ‘number of wheels’ is compatible with both trees.



**Figure 3.1:** Putative phylogenies of vehicles. Tree (a) is compatible with the character ‘existence of engine’. The engine is invented once, in the edge connecting the root and the common ancestor of car and motorcycle. Tree (b) is compatible with ‘number of wheels’, but would be incompatible with ‘existence of engine’.

**Exercise:** The objects  $A, B, C, D$  share three characters 1,2,3. The following matrix holds their states:

	1	2	3
A	$a$	$\alpha$	$d$
B	$a$	$\beta$	$e$
C	$b$	$\beta$	$f$
D	$b$	$\gamma$	$d$

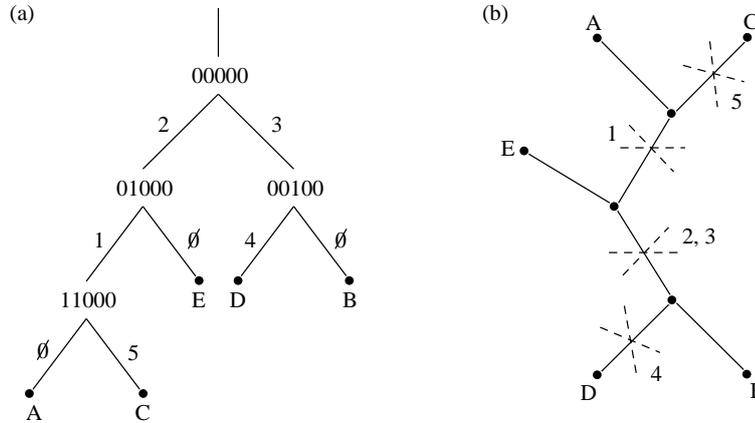
Look at all possible tree topologies. Is there, among all these trees, a tree  $T$  such that all characters are compatible with  $T$ ?

### 3.3 Perfect Phylogenies

Let a set  $\mathcal{C}$  of characters and a set  $S$  of objects be given.

**Definition:** A tree  $T$  is called a *perfect phylogeny (PP)* for  $\mathcal{C}$  if all characters  $C \in \mathcal{C}$  are compatible with  $T$ .

**Example:** The objects  $A, B, C, D, E$  share five binary (two-state) characters. The matrix  $M$  holds their binary states:



**Figure 3.2:** Perfect phylogeny for binary character states of  $M$  shown as (a) a rooted tree (numbers at edges denote invention of characters), and (b) an unrooted tree (crossed dashed lines represent splits, a number at a split denotes the character inducing the split)

	1	2	3	4	5
A	1	1	0	0	0
B	0	0	1	0	0
C	1	1	0	0	1
D	0	0	1	1	0
E	0	1	0	0	0

The trees in Figure 3.2 show a perfect phylogeny for  $M$ . There are two ways of looking at this perfect phylogeny:

- As a rooted tree in a directed framework with development (Figure 3.2 (a)): A (virtual) ancestral object is added as the root. Its characters have state 00000. During development, on the way from this root to the leaves, a character may invent a new state '1' exactly once. In the subtree below the point where the state switches, the state will always remain '1'.
- As an unrooted tree in a split framework (Figure 3.2 (b)): A binary character corresponds to a split (a bipartition) in a set of objects. In a tree this is represented by an edge. In this framework, characters 4 and 5 are *uninformative*. The splits they induce are trivial as they separate a single object from the other ones and hence do not tell us about the relationships among the objects. Character 3 is just the opposite of character 2. It describes the same split.

We can define an *object set* for each character: character  $i$  corresponds to the set of objects  $O_i$  where the character  $i$  is "on", that is it has value 1. Then each column of the matrix  $M$  corre-

sponds to a set of objects:  $O_1 = \{A, C\}$ ,  $O_2 = \{A, C, E\}$ ,  $O_3 = \{B, D\}$ ,  $O_4 = \{D\}$ ,  $O_5 = \{C\}$ .

The *Perfect Phylogeny Problem (PPP)* addresses the question if for a given matrix  $M$  there exists a tree as shown in Figure 3.2 (a) and, given its existence, how to construct it.

An alternative formulation for the PPP with binary characters is the formulation as a *Character Compatibility Problem*: Given a finite set  $S$  and a set of splits (binary characters)  $\{A_i, B_i\}$  such that  $A_i \cap B_i = \emptyset$  and  $A_i \cup B_i = S$ , is there a tree as in Figure 3.2 (b) that realizes these splits?

In general, the PPP is NP-hard. Given binary characters, it is easy, though. Gusfield [8] formulates the following theorem:

**Theorem.** If all characters are binary,  $M$  has a PP if and only if for any two columns  $i, j$  there holds one of:

- (i)  $O_i \subseteq O_j$ ,
- (ii)  $O_j \subseteq O_i$ ,
- (iii)  $O_i \cap O_j = \emptyset$ .

The condition in the above theorem describes the compatibility of two characters. Two characters are compatible if they allow for a PP. And Gusfield's theorem becomes: "A set of binary characters has a PP if and only if all pairs of characters are compatible." Recently, Lam *et al.* [9] showed that this sentence similarly holds for three-state characters: "There is a three-state PP for a set of input sequences if and only if there is a PP for every subset of three characters." But note that this is not true in general, i.e. for  $r$ -state characters with  $r > 3$ , as there are many "tree realizations" (see [10]).

The naïve algorithm to check if a PP exists is to test all pairs of columns for the above conditions which has time complexity  $O(nm^2)$ , where  $n$  is the number of rows, and  $m$  is the number of columns in  $M$ . We will present here a more sophisticated method for recognition *and* construction of a PP with a time complexity of  $O(nm)$ . The algorithm is influenced by Gusfield ([8] and [3, Chapter 17.3.4]), Waterman [11], and Setubal and Meidanis [12].

The algorithm is as follows (an example is given on page 20):

1. Check the inclusion/disjointness property of  $M$ :
  - a) Sort the columns of  $M$  by their number of ones.  
(Note that this can also be done in  $O(nm)$  time.)
  - b) Add a column number 0 containing only ones.

- c) For each '1' in in the sorted matrix, draw a pointer to the previous '1' in the same row.

**Observation.**  $M$  has a PP if and only if in each column, all pointers end in the same preceding column.<sup>1</sup>

If the condition is fulfilled, we can proceed constructing a PP.

2. Build the phylogenetic tree:

- a) Create a graph with a node for each character.
- b) Add an edge  $(i, j)$  between direct successors in the partial order defined by the  $O_i$  set inclusion as indicated by the pointers: If the pointers of column  $j$  end in column  $i < j$ , draw an edge from  $i$  to  $j$ . Because of the set inclusion/disjointness conditions, the obtained graph will form a tree.

3. Refine and annotate the tree:

- a) Annotate all leaf nodes labeled  $j$  with the object set  $O_j$ , mark these objects as 'used', and annotate the parent edge with  $j$ .
- b) In a bottom-up fashion, for each inner node which is labeled with character  $j$  and whose children are all re-labeled, do
- i. Re-label the node with all objects of  $O_j$  which are not yet marked as 'used',
  - ii. mark these objects as 'used', and
  - iii. annotate the parent edge with  $j$  (not for the root node).
- c) For each node  $u$  and each object  $o$  it is labeled with, if  $u$  is not a leaf labeled only with  $o$ , remove  $o$  from the labeling and append a new leaf  $v$  labeled with  $o$ .

Finally, we obtain a PP where the edges are annotated with the invented characters and the leaves are labeled with the objects.

<sup>1</sup>Thanks to Zsuzsanna for suggesting the pointer trick.

### 3 Characters, States, and Perfect Phylogenies

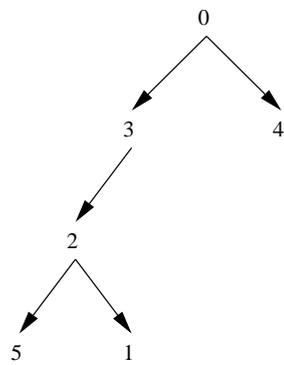
#### Example:

	1	2	3	4	5
A	0	1	1	0	1
B	0	0	0	1	0
C	1	1	1	0	0
D	0	0	0	1	0
E	0	0	1	0	0

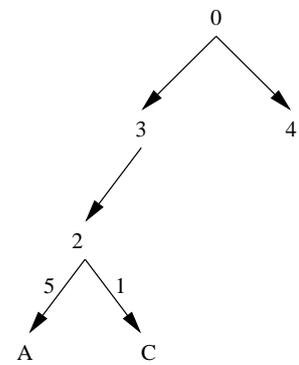
$\xRightarrow{1.}$

	0	3	2	4	1	5					
A	1	←	1	←	1	←	0	←	0	←	1
B	1	←	0	←	0	←	1	←	0	←	0
C	1	←	1	←	1	←	0	←	1	←	0
D	1	←	0	←	0	←	1	←	0	←	0
E	1	←	1	←	0	←	0	←	0	←	0
# '1's		3	2	2	1	1					

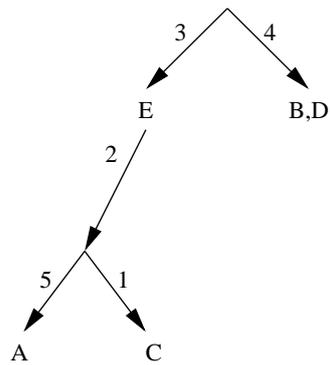
$\xRightarrow{2.}$



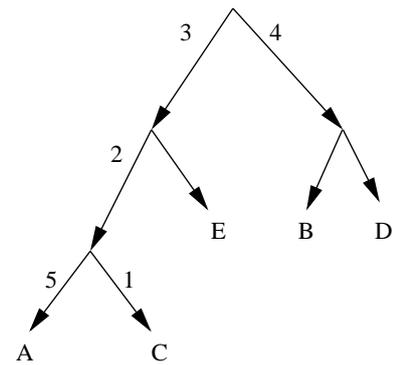
$\xRightarrow{3.a)}$



$\xRightarrow{3.b)}$



$\xRightarrow{3.c)}$



**Part II**

**Parsimony**

# 4 The Small Parsimony Problem

## 4.1 Introduction

There are many possibilities to construct a phylogenetic tree from a set of objects. Of course, we want to find the ‘best’ tree, or at least a good one. Judging the quality of a phylogenetic tree requires criteria. *Parsimony* is such a criterion.

The general idea is to find the tree(s) with the minimum amount of evolution, i.e., with the fewest number of evolutionary events. More generally, we use some weighting scheme to assign specific costs to each event and seek for an evolutionary scenario with minimal total cost. We call such a tree the *most parsimonious tree*.

There are several methods to reconstruct the most parsimonious tree from a set of data. First, they have to find a possible tree. Second, they have to be able to calculate and optimize the changes of state that are needed. As the second problem is the easier subproblem of the general parsimony problem, we start with this one: We assume that we are given a tree topology and character states at the leaves, and the problem is to find a most parsimonious assignment of character states to internal nodes. This problem is referred to as the *small parsimony problem*.

Not only the concept of parsimony itself, but a considerable amount of variants of weighting schemes for different applications has been studied and discussed in the past. Traditionally, the following properties are used to characterize a cost function  $cost(c, e, x, y)$ , which assigns a specific cost to the event “character  $c$  changes at edge  $e$  from state  $x$  to  $y$ ”.

**Character dependency.** A parsimony cost function  $cost$  is called *character independent* if  $cost(c_i, e, x, y) = cost(c_j, e, x, y)$  for any characters  $c_i$  and  $c_j$ .

**Edge dependency.** It is called *edge independent* if  $cost(c, e_i, x, y) = cost(c, e_j, x, y)$  for any edges  $e_i$  and  $e_j$ . If the evolutionary distances in a tree are known, e.g. in terms of time and/or mutation rate, this knowledge can be used to define costs for events on each edge separately.

**Symmetry.** It is called *symmetric* if  $cost(c, e, x, y) = cost(c, e, y, x)$  for all states  $x$  and  $y$ . We might assume that the invention of a phenotype ( $0 \rightarrow 1$ ) is less probable than its

loss ( $1 \rightarrow 0$ ). This can easily be modeled by assigning a higher value to  $cost(c, e, 0, 1)$  than to  $cost(c, e, 1, 0)$ .

In the following, we want to shortly review some classical, concrete weighting schemes. For a broader and more detailed discussion, see for example the review of Felsenstein [5].

In 1965, Camin and Sokal introduced a model that assumes an ordering of the states. A character can only change its state stepwise in the given direction—from “primitive to derived”—and reversals are not allowed. Another long-standing model is the so-called *Dollo Parsimony*. It assumes that the state 1 corresponds to a complex phenotype. Following Dollo’s principle, the loss of such a rare property is irreversible. Thus *homoplasy*, the development of a trait in different lineages, is precluded. That means, state 1 is allowed to arise only once and we minimize only the number of changes from 1 to 0. Corresponding theory and methods have been developed by Farris and Le Quesne in the 1970s. However, there are several counterexamples for homoplasy known by now, for example the parallel development of wings of birds and bats. Since we cannot exclude this phenomenon in the development of phenotypes, we should not rule it out on principle, but rather explicitly allow homoplasy to take place. This ensures to be able to detect and examine this interesting feature, which might elucidate the evolution of certain phenotypes.

The most basic but also commonly used cost function for multiple state characters is the *Fitch parsimony* (Felsenstein defines this model as *Wagner parsimony*). Here, we simply count the number of events by assuming (character and edge independent) *unit costs*: If a character has state  $x$  at some node and state  $y$  at a child node, we assign  $cost(x, y) = 0$  for  $x = y$  and  $cost(x, y) = 1$  for  $x \neq y$ .

A common method for finding a most parsimonious labeling under the Fitch parsimony weighting scheme was first published by Fitch and is thus well-known as the *Fitch algorithm*. Concurrently, Hartigan developed a similar, but more powerful framework, published two years later in 1973. In contrast to the Fitch algorithm, his method is not restricted to binary trees, and can be used to find all co-optimal solutions. Since he proved its correctness, he also first showed the correctness of Fitch’s algorithm by presenting it as a special case of his general method.

Shortly after, Sankoff and Rousseau described a generalization for any edge independent, arbitrary metric cost function which may vary for certain events. In this context, where not only the number of changes is counted, but the “amount of change” is weighted, we also speak of the *minimal mutation problem*. Erdősch and Székely, in turn, proposed a similar approach in a further generalized framework for edge dependent weights in 1994. Even for asymmetric models, a linear-time algorithm was introduced by Miklós Csűrös recently.

All of the above variants of the small parsimony problem can be subsumed in a formal problem statement: We have given a phylogenetic tree  $T = (V, E)$  where each leaf  $l$  is labeled

with state  $s(c, l)$  for character  $c$ . We want to assign a state  $s(c, v)$  to each inner node  $v$  for each character  $c$  such that the total cost  $W(T, s)$  is minimized:

$$W(T, s) := \sum_{(u,v) \in E} \sum_{c \in \mathcal{C}} \text{cost}(c, (u, v), s(c, u), s(c, v))$$

In general, all models assume that the different characters evolve independently. This allows us to perform the construction of a parsimonious labeling character-wise. That means, we can compute labelings for all characters separately, each minimizing the parsimony cost. Finally, the combination results in an overall labeling which minimizes the total cost.

$$\begin{aligned} \min_s(W(T, S)) &= \min_s \left( \sum_{(u,v) \in E} \sum_{c \in \mathcal{C}} \text{cost}(c, (u, v), s(c, u), s(c, v)) \right) \\ &= \sum_{c \in \mathcal{C}} \min_s \left( \sum_{(u,v) \in E} \text{cost}(c, (u, v), s(c, u), s(c, v)) \right) \end{aligned}$$

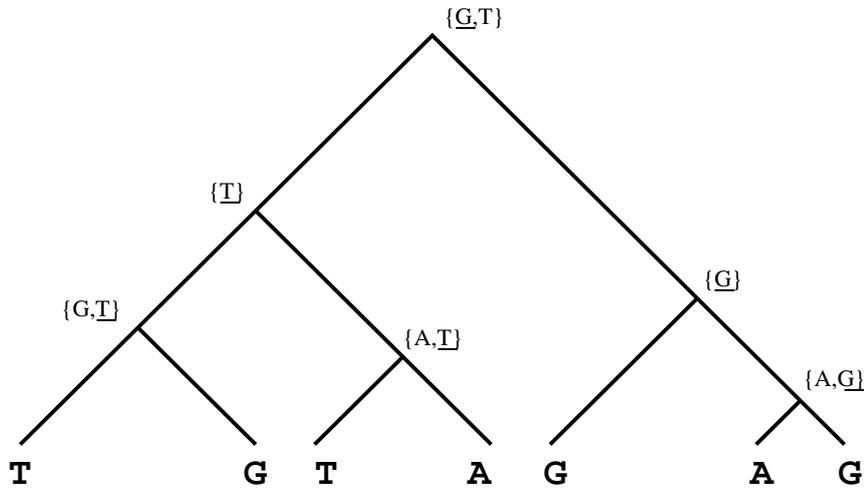
In the following, we present three classical algorithms: We begin with the simple algorithm of Fitch [13] and then proceed with the generalized method of Hartigan [14]. Finally, we introduce a dynamic programming algorithm due to Sankoff and Rousseau [15] (nicely described in [16]) that works for more general state change costs.

## 4.2 The Fitch Algorithm

Fitch's algorithm solves the small parsimony problem for unit costs (wagner parsimony) on binary trees. Recall the problem statement: We have given a phylogenetic (binary) tree  $T = (V, E)$  where each leaf  $l$  is labeled with state  $s(c, l)$  for character  $c$ . We want to assign a state  $s(c, v)$  to each inner node  $v$  for each character  $c$  such that the total number of state changes is minimal.

We assume that the tree is rooted. If not, we can introduce a root node at any edge without altering the result. The method works in a dynamic programming fashion. For each character, it performs two main steps to construct a labeling  $s(c, v)$  for each node  $v$  of the tree: A bottom-up phase, followed by a top-down refinement.

In the following we give the algorithm for processing one character. It is performed for each character separately to obtain an optimal overall labeling. An example can be found in Figure 4.1 on the next page.



**Figure 4.1:** Example for the Fitch algorithm on one character. For each internal node, the candidate set is given. The bars indicate one of two possible solutions which can be found during the top-down refinement.

**1. (Bottom-up phase)** We traverse the tree from the leaves to the root such that when a node is processed, all its children have already been processed. Obviously the root is the last node processed by this traversal. During this phase, we collect putative states for the labeling of each node  $v$ , stored in a candidate set  $S(v)$ .

**1a. (Leaves)** For each leaf  $l$ , set  $S(l) := \{s(c, l)\}$ .

**1b. (Internal nodes)** Assume an internal node  $u$  with children  $v$  and  $w$ . If the  $v$  and  $w$  share common candidates, these are candidates for  $u$  as well. Otherwise, the candidates of both children have to be considered as candidates for  $u$ . More formally:

$$S(u) := \begin{cases} S(v) \cap S(w) & \text{if } S(v) \cap S(w) \neq \emptyset, \\ S(v) \cup S(w) & \text{otherwise.} \end{cases}$$

**2. (Top-down refinement)** The most parsimonious reconstruction of states at the internal nodes is then obtained in a top-down pass according to the following rules:

**2a. (Root)** If the candidate set of the root contains more than one element, arbitrarily assign one of these states to the root.

**2b. (Other nodes)** Let  $v$  be a child of node  $u$ , and let  $a$  denote the state assigned to  $u$ .

(i) If  $a$  is contained in  $S(v)$ , assign it to node  $v$  as well.

(ii) Otherwise, arbitrarily assign any state from  $S(v)$  to node  $v$ .

Finally, a most parsimonious overall labeling  $s$  can be returned. A traversal takes  $O(n)$  steps and the computation for each node can be done in time  $O(\sigma)$ , where  $\sigma$  is the maximal number of different states for a character. Hence, the two phases can be performed in  $O(n\sigma)$  time for each of the  $m$  characters. This gives a linear total time complexity of  $O(mn\sigma)$ .

This kind of two step approach is generally called *backtracing*: In a forward phase, intermediate results are computed, which are then used (like a “trace”) during a backward phase to compose a final solution. In steps 2a and 2b of the above algorithm, we generally can choose between different states to label a node. In any case, we obtain an optimal solution. But, for this algorithm, *backtracking* (recursively following each possible way of constructing a solution) does *not* yield all optima.

### 4.3 Hartigan’s Algorithm

In contrast to the Fitch algorithm, Hartigan’s algorithm is not restricted to binary trees and backtracking would yield *all* optima. Again, we assume unit costs. A careful comparison to Fitch’s method will expose that it is as a special case of Hartigan’s.

Here, we only give a motivation for how and why the algorithm computes an optimum. For a more detailed discussion and the proof of correctness, we refer to the original publication of Hartigan [14].

Similar to Fitch, for each character, the algorithm computes candidate states in a bottom-up phase and determines an optimum in a top-down refinement. See Figure 4.2 on the facing page for an example. Analogously, the running time is in  $O(mn\sigma)$ .

**1. (Bottom-up phase)** During a traversal beginning at the leaves and going up to the root, we collect candidate states for the labeling of each node, where putative states are stored in two sets: Analog to Fitch, candidate states which definitely yield a most parsimonious labeling are collected in a set  $S_1$ , the primary set. But additionally, a secondary set  $S_2$  is used, to deposit states which alternatively could also result in a correct solution.

**1a. (Leaves)** For a leaf  $l$ , the state is given and therefore fixed. Hence the primary candidate set is  $S_1(l) := \{s(c, l)\}$ , and there is no alternative labeling,  $S_2(l) := \emptyset$ .

**1b. (Internal nodes)** To determine the candidate sets for an internal node  $u$ , we first count the occurrences of each state in the descendent nodes. On the one hand, the selection of the majority will imply the fewest events and thus give a valid solution; even if finally the label of its parent node differs from this state, which adds an unavoidable extra weight of one. But on the other hand, in the latter case, choosing another state might not increase the weight at the parent edge if

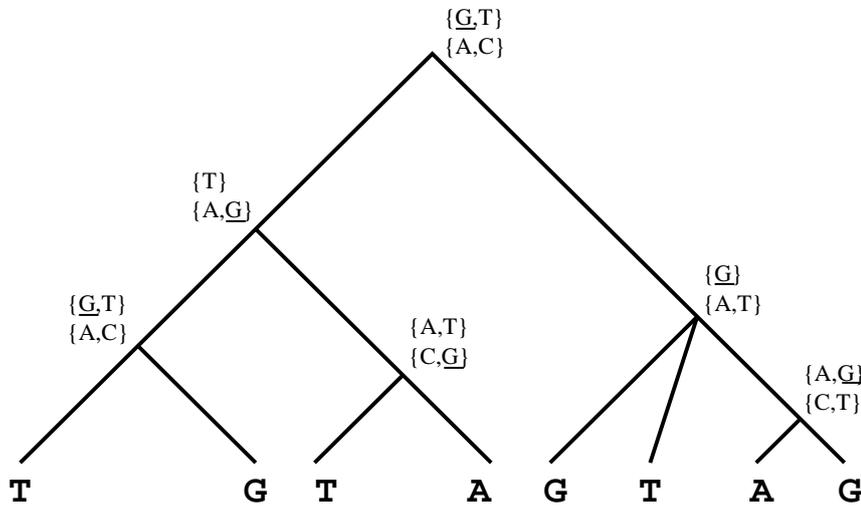
the labels are equal. If, additionally, this other state occurs exactly once fewer in the child nodes than the maximum, this would also yield a total increase by one. To provide the opportunity to find *all* solutions, we store the majority state (or all such states in the case of a tie) in  $S_1(u)$ , and the putative alternatives in  $S_2(u)$ . Formally, we proceed as follows:

$$\begin{aligned} \text{For each state } b, \text{ define } k(b) &:= |\{v_i \mid b \in S_1(v_i)\}| \\ K &:= \max_b \{k(b)\} \\ S_1(u) &:= \{b \mid k(b) = K\} \\ S_2(u) &:= \{b \mid k(b) = K - 1\} \end{aligned}$$

**2. (Top-down refinement)** In this second phase, we have to decide for each node, which of its candidate states we select to obtain a most parsimonious labeling.

**2a. (Root)** For the root node  $r$ , selecting any of the states from the primary set  $S_1(r)$  will finally yield a proper solution if the descendent nodes are labeled correspondingly.

**2b. (Other nodes)** Then, in a top-down traversal, we assign a labeling to each node  $v$  depending on the labeling of its parent node  $u$ . If the ancestral state  $a$  is also a candidate in the primary set  $S_1(v)$ , only this state can minimize the number of events. But if  $a$  is contained in the secondary set  $S_2(v)$ , both  $a$  and a state in



**Figure 4.2:** An example for Hartigan's algorithm on one character. For each internal node, the primary and secondary set is given. The bars indicate one possible solution. Further optima exist.

$S_1(v)$  imply exactly one unavoidable state change: Either one extra change on one of the descendent edges, or a change on the edge  $(u, v)$ . Otherwise, if the ancestral state  $a$  is not contained in any of the two candidate sets, selecting  $a$  for  $v$  would imply at least two changes among  $v$  and its descendants more than the state in  $S_1(v)$ . Hence in this case, the most parsimonious choice is a state in  $S_1(v)$ , even though this causes a change from  $s(c, u)$  to  $s(c, v)$ .

Hence, we refine  $v$  to  $s(c, v) := b$ , with any  $b \in B(v, a)$  where

$$B(v, a) := \begin{cases} \{a\} & \text{if } a \in S_1(v), \\ \{a\} \cup S_1(v) & \text{if } a \in S_2(v), \\ S_1(v) & \text{otherwise.} \end{cases}$$

## 4.4 The Sankoff Dynamic Programming Algorithm

The above algorithms are restricted to unit costs. Sankoff developed a generalization for symmetric, edge independent costs which may vary for certain events [17, 16]. Again, a backtracing approach, composed of a bottom-up and a top-down phase, is performed for each character separately. Similar to Hartigan's algorithm, backtracking can be used to obtain all optimal solutions.

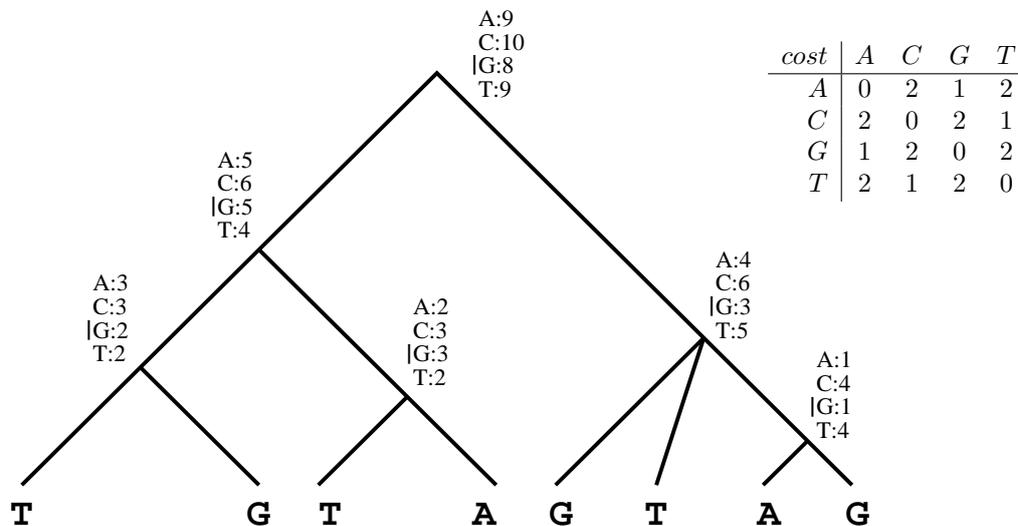
The algorithm is as follows (see Figure 4.3 on the next page for an example).

**1. (Bottom-up phase)** The tree is traversed bottom-up. During this traversal, assume we process a node  $u$ . Define  $C(u)$  as the cost of the optimal solution of the minimum mutation problem for the subtree  $T_u$  rooted at  $u$ . Let  $C(u, a)$  be the cost of the best labeling of  $T_u$  when node  $u$  is required to be labeled with state  $a$ . Obviously,  $C(u) = \min_a C(u, a)$ .

**1a. (Leaves)** The state for a leaf  $l$  is fixed by the input. For this base case we set  $C(l, a) = 0$  for  $a = s(c, l)$  and  $C(l, a) = \infty$ , otherwise.

**1b. (Internal nodes)** The recurrence relation to compute  $C(u, a)$  for an internal node  $u$  is given by

$$C(u, a) = \sum_{v \text{ child of } u} \min_b (\text{cost}(a, b) + C(v, b)).$$



**Figure 4.3:** The dynamic programming algorithm of Sankoff on one character for the given weighting scheme. For each internal node, the values for  $C$  are shown. The bars indicate a solution. In this case, there is no further optimum.

**2. (Top-down refinement)** The optimal assignment of states to the internal nodes is then obtained in a backtracing phase.

**2a. (Root)** The root  $r$  is assigned a state  $a$  such that  $C(r) = C(r, a)$ .

**2b. (Other nodes)** In a top-down traversal, the child  $v$  of an already labeled node  $u$  (say,  $u$  was labeled with state  $a$ ) is assigned a state  $b$  that yielded the minimum in the bottom-up pass, i.e., where

$$cost(a, b) + C(v, b) = \min_{b'} [cost(a, b') + C(v, b')].$$

In contrast to the previous algorithms, processing one node in the bottom-up phase (step 1b.) takes  $O(\sigma^2)$  steps: For each state in  $u$  we have to minimize over all possible states in  $v$ . Hence, the algorithm solves the problem correctly in  $O(mn\sigma^2)$  time and space, where  $m$  is the number of characters,  $n$  the number of leaves, and  $\sigma$  is the maximal number of different states of a character.

# 5 Maximum Parsimony

## 5.1 Introduction

The algorithms covered in the previous chapter solve the “small parsimony problem”, where a tree is given. They compute the most parsimonious assignment of states to the inner nodes of the tree. However, in reality we do not know the tree, but are looking for the tree that yields the minimum cost (maximum parsimony) solution among all possible tree topologies.

Therefore, the algorithms that solve the *Maximum Parsimony Problem* find the optimal tree among all possible ones. The problem is actually NP-hard. Nevertheless, we present some slow and exact methods and some faster but possibly suboptimal heuristics that solve the problem. Finally, we mention an even harder problem (generalized tree alignment) that arises when no alignment is given initially. We also discuss drawbacks of the parsimony model (Section 5.5).

Before describing algorithms, however, we clarify the meaning of the term *heuristic*. In general, one speaks of a heuristic method for the solution of a problem if the method uses exploratory problem-solving techniques based on experience or trial-and-error, without in the worst case necessarily to perform better or faster than naive methods. It can mean two different things.

**Correctness heuristic.** A fast (efficient) computational method for a computationally difficult problem that is not guaranteed to return an optimal solution, but is expected to produce a “good” solution, sometimes provably within a certain factor of the optimum.

**Running time heuristic.** An exact algorithm that is guaranteed to return the optimal solution, but not within a short time, although it may do so on some inputs.

An example of the first type of heuristic is the greedy method discussed in Section 5.2.3, and an example of the second type of heuristic is the branch-and-bound method discussed in Section 5.2.2.

$m$  columns

A	...	Q	N	L	A	K	R	G	H	N	N	Y	K	...
B	...	H	K	L	A	-	-	V	A	N	N	I	K	...
C	...	E	N	M	A	K	R	G	R	N	N	Y	N	...
D	...	A	E	M	A	E	R	G	-	-	N	L	A	...

**Figure 5.1:** A multiple alignment for five taxa A, B, C, D with  $m$  columns.

## 5.2 Exact Methods

### 5.2.1 Brute-Force Method: Exhaustive Enumeration

The algorithms for the small parsimony problem allow us to evaluate any given tree topology. Therefore a naive algorithm suggests itself: In principle we can solve the small parsimony problem for *all possible* tree topologies, and the one(s) with lowest cost will be the most parsimonious tree(s). As we have already seen, the number of different trees grows super-exponentially with the number of taxa; therefore an exhaustive enumeration is in practice infeasible for more than about 12 taxa. However, the advantage would be that we get all most parsimonious trees if there are several equally good solutions.

### 5.2.2 Branch and Bound Heuristic.

One speaks of a *branch-and-bound* heuristic whenever it is possible to restrict the search such that not the complete theoretically possible search space has to be examined, but one can already stop whenever a partial solution provably can not lead to an (optimal) solution of the complete problem. The maximum parsimony problem is well suited to apply the branch-and-bound heuristic.

If we construct the tree from a multiple alignment where we consider each column as an independent character, there are two alternative ways to apply branch and bound heuristics to the problem of finding the most parsimonious tree. Assume the alignment has  $m$  columns, as shown in Figure 5.1. The first strategy does branch and bound on the alignment columns, and the second one does branch and bound on the alignment rows, i.e. the leaves of the tree. Both procedures have in common that they start with an upper bound on the length of the most parsimonious tree, obtained e.g. by an approximation algorithm such as a minimum spanning tree (see Section 5.3.3).

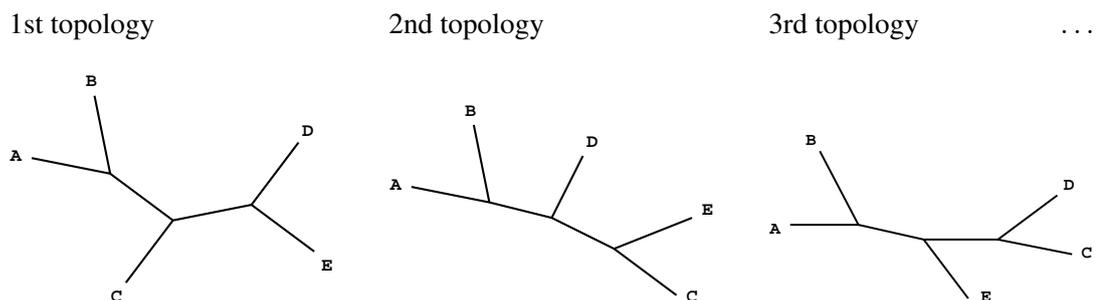
### Column-wise Branch-and-Bound

Intuitively it seems reasonable to apply branch and bound to the columns of the sequence alignment, since they are independent by assumption. For each possible tree topology, one would first compute the minimal length for the first alignment column, then for the second, and so on. Once the sum of these lengths is larger than the best solution computed so far, the procedure stops and continues with the next tree topology (see Figure 5.2). The speedup by this procedure is not impressive, though.

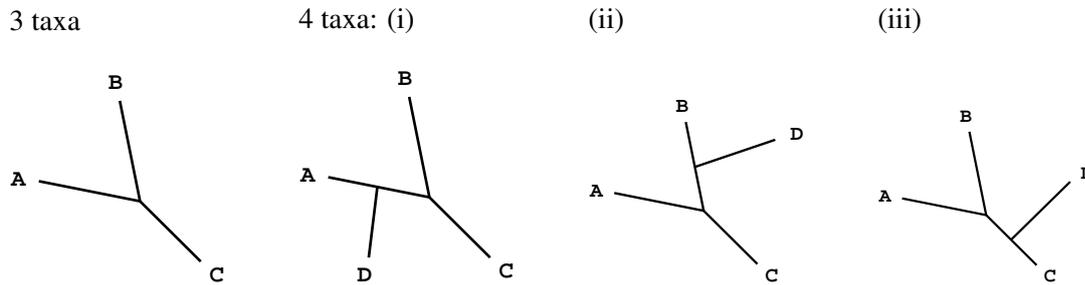
### Row-wise Branch-and-Bound

The second approach, which is due to Hendy and Penny [18], does branch and bound on the sequences, in parallel to the enumeration of all possible tree topologies. The idea is that adding branches to a tree can only increase its length. Hence we start with the first three taxa, build the only possible unrooted tree for them and compute its length. Then we add the next taxon in the three possible ways, thereby generating the three possible unrooted trees for four taxa (see Figure 5.3). Whenever one of the trees already has length larger than the best solution computed so far, the procedure stops, otherwise it is continued by adding the next taxon, etc. This procedure may be used for finding the maximum parsimony tree for up to 20 taxa. The actual speedup that can be achieved depends on the input alignment, however. In a worst case scenario, the actual speedup may be quite small.

For even larger numbers of taxa, heuristic procedures have to be used that no longer guarantee optimality, some ideas can be found in the following section.



**Figure 5.2:** Branch and bound applied to the alignment columns.



**Figure 5.3:** Branch and bound applied to the alignment rows.

### 5.2.3 Greedy Sequential Addition

Adding a taxon in all possible edges of an existing tree, thereby generating a number of new trees, is called *sequential addition*. This approach is used to enumerate the tree space in the row-wise Branch-and-Bound scenario in Section 5.2.2. The idea can be used in an even faster, but less accurate greedy approach: Starting with a tree with three taxa, a fourth taxon is added to each possible edge, and the net amount of evolutionary change is determined (e.g., by using the Fitch algorithm). However, only the best of the resulting trees is kept. It is used as input in the next iteration.

Finally, once the (heuristic) tree is constructed, it can be improved in an iterated fashion by local modifications including *branch swapping* and *subtree pruning and regrafting*.

The best tree computed in this way is not necessarily the most parsimonious one, but mostly a good approximation. The quality of the resulting tree also depends on the choice of the three taxa that are used for the initial tree [19].

## 5.3 Steiner Trees and Spanning Trees

Steiner trees can also be used to solve the Maximum Parsimony Problem. However, finding the best Steiner tree is itself a difficult problem. Before we explain it (see Section 5.3.6), we introduce some basic concepts.

### 5.3.1 The DNA Grid Graph

If the objects whose phylogenetic tree we wish to reconstruct are aligned DNA sequences, all of the same length  $n$ , the *DNA grid graph* provides the means to formalize the notion of such a tree “as perfect as possible”.

**Definition.** The *DNA grid graph*  $G = (V, E)$  is defined as follows. Each DNA sequence of length  $n$  is represented by a vertex in  $V$ . Two vertices are connected by an edge in  $E$  if and only if there is exactly one mismatch between the represented sequences.

Obviously, the graph contains  $4^n$  vertices, and the length of a shortest path between two vertices is the *Hamming distance* between the represented sequences.

Now suppose that we are given  $N$  aligned sequences of length  $n$  each. These define the set  $U$  of *terminal nodes* in this graph.

### 5.3.2 Steiner Trees

**Definition.** Given a connected weighted graph  $G = (V, E)$  and a subset of the vertices  $U \subseteq V$  of terminal nodes. A tree that is a subgraph of  $G$  and connects the vertices in  $U$  is called a *Steiner tree* of  $U$ .

Note that in general, the Steiner tree will also contain a subset of the remaining vertices of  $G$ , not in  $U$ , which we call *Steiner nodes*.

**Definition.** A Steiner tree of minimal length is a *minimum Steiner tree*.

**Minimum Steiner Tree Problem.** Given a connected weighted graph  $G = (V, E)$  and a subset of the vertices  $U \subseteq V$ , find a minimum Steiner tree of  $U$ .

**Observation.** Assume once again that the set of terminal vertices  $U$  are objects from the DNA grid graph. Then, the minimum Steiner tree of  $U$  will be a shortest tree that connects all vertices in  $U$ , and hence, explains the DNA sequence data.

**Theorem [20].** The Minimum Steiner Tree Problem is NP complete.

### 5.3.3 Spanning Trees

**Definition.** Let  $G = (V, E)$  be a connected weighted graph. A tree that is a subgraph of  $G$  and connects all vertices of  $G$  is called a *spanning tree*. A spanning tree of minimum length is called a *minimum spanning tree*.

**Property.** Let  $T$  be a minimum spanning tree in a graph  $G$ . Let  $e$  be an edge in  $T$ , splitting  $T$  into two subtrees  $T_1$  and  $T_2$ . Then  $e$  is of least weight among all the edges that connect a node of  $T_1$  and a node of  $T_2$ .

There are fairly simple algorithms to construct a minimal spanning tree. Two such algorithms are presented in Section 24.2 of [21], Kruskal's algorithm and Prim's algorithm. Both algorithms run in  $O(|E| \log |V|)$  time. They are based on the following *generic algorithm*:

Given a graph  $G$ , maintain a set of edges  $A$  that in the end will form the minimum spanning tree. Initially,  $A$  is empty. Then step by step *safe* edges are added to  $A$ , in the sense that an edge  $e = \{u, v\}$  is safe for  $A$  if  $A \cup \{e\}$  is a subset of a minimum spanning tree.

**Kruskal's Algorithm** implements the selection of a safe edge as follows: Find, among all edges that connect two trees in the growing forest that one of least weight.

Correctness: Follows from the above property of minimum spanning trees.

Implementation:

1. Sort the edges in  $E$  by nondecreasing weight, such that  $w(e_1) \leq w(e_2) \leq \dots \leq w(e_{|E|})$ .
2. For each edge  $i$  from 1 to  $|E|$ , test if  $e_i$  forms a circle (for example by maintaining for each node a representative element from the connected component that it is contained in), and if not, add it to  $A$ .

The first step takes  $O(|E| \log |E|)$  time. For the second step, a run time of  $O(|E| \alpha(|E|, |V|))$  can be achieved using a disjoint-set data structure where  $\alpha$  is the (very slowly growing) inverse of Ackermann's function. For details, see [21].

**Prim's Algorithm** also follows the generic algorithm, using a simple greedy strategy. The selected edges always form a single tree. The tree starts from an arbitrary root vertex and at each step an edge of minimum weight connecting a vertex from the tree with a non-tree vertex is selected. This is repeated until the tree spans all vertices in  $V$ .

Correctness: Follows from the above property of minimum spanning trees.

Implementation: Maintain a priority queue that contains all vertices that are not yet members of the tree, based on the least weight edge that connects a vertex to the tree. Using a binary heap, the priority queue can be implemented in  $O(|V| \log |V|)$  time, resulting in a total running time of  $O(|E| \log |V|)$ . (This can be improved to  $O(|E| + |V| \log |V|)$  by using a Fibonacci heap.) Again, see [21] for details.

### 5.3.4 Spanning Trees and the Traveling Salesman Problem

**Definition.** Given a connected graph  $G$ , a simple cycle that traverses all the nodes of  $G$  is called a *Hamiltonian cycle*.

In a connected weighted graph, a Hamiltonian cycle of minimum length is also called a *Traveling Salesman Tour*.

**The Traveling Salesman Problem (TSP).** Given a connected graph  $G$ , find a Traveling Salesman Tour.

**Theorem.** The Traveling Salesman Problem is NP hard. (One of the classic NP-hard problems.)

A simple approximation algorithm for the Traveling Salesman Problem on a metric is the following one:

1. Compute a minimum spanning tree.
2. Traverse the minimum spanning tree in a circular order; each time a node is visited for the first time, output it.

**Definition.** A *constant factor approximation algorithm* is a heuristic algorithm for the solution of a problem such that the result of the heuristic will deviate from the optimal solution by at most a certain multiplicative factor.

For example, a 2-approximation of a minimizing optimisation problem is an algorithm whose solution is guaranteed to be at most twice the optimal solution. An example is the following:

**Theorem.** The above tour gives a 2-approximation for the Traveling Salesman Problem.

**Proof.** Let  $H^*$  be an optimal tour,  $T$  be a minimum spanning tree,  $W$  be a full walk of  $T$ , and  $H$  be the tour output by the above algorithm. Then we have:  $c(T) \leq c(H^*)$  because by removing one edge from  $H^*$  one gets a (linear) tree which is not shorter than the MST. Obviously,  $c(H) \leq c(W) = 2c(T)$ , and hence  $c(H) \leq 2c(H^*)$ .  $\square$

### 5.3.5 Spanning Trees as Approximations of Steiner Trees

A simple approximation algorithm for the Minimum Steiner Tree Problem on a set of vertices  $U$  is the *spanning tree heuristic*:

1. Compute all-pairs shortest-paths on the terminal nodes in  $U$ . This takes  $O(|V|(|V| + |E|))$  time.
2. Compute a minimum spanning tree on this complete, weighted subgraph. As shown above, this takes  $O(|U|^2 \log |U|)$  time.
3. Map back the minimum spanning tree into the original graph. This step takes  $O(|E|)$  time.

**Theorem.** The length of the resulting Steiner tree is at most twice the length of the minimum Steiner tree.

**Proof.** Let  $T$  be a minimum spanning tree and  $T^*$  be a minimum Steiner tree. Let  $W^*$  be a full walk of  $T^*$ . We want to show that  $c(T) \leq 2c(T^*)$ . Let  $H^*$  be a Travelling Salesman Tour on  $U$  and the Steiner nodes of  $T^*$ . We make the following two observations: (1)  $c(T) \leq c(H^*)$  (see the proof of the above theorem) and (2)  $2c(T^*) = c(W^*) \geq c(H^*)$  (this follows from the definition of the Travelling Salesman tour). Together we have  $c(T) \leq c(H^*) \leq 2c(T^*)$ .  $\square$

### 5.3.6 Application to Phylogeny

**Summary:** Given a set of aligned DNA sequences whose phylogeny we want to reconstruct, we represent them as vertices in a DNA grid graph. Now we can use the spanning tree heuristic to compute a Steiner tree on these vertices. This tree is a 2-approximation of the correct Steiner tree, i.e., one with the minimum amount of evolutionary changes. BUT: The approximated Steiner tree is not necessarily binary, nor will it necessarily have the terminal nodes as leaves (which would be expected of a phylogenetic tree).

## 5.4 Generalized Tree Alignment

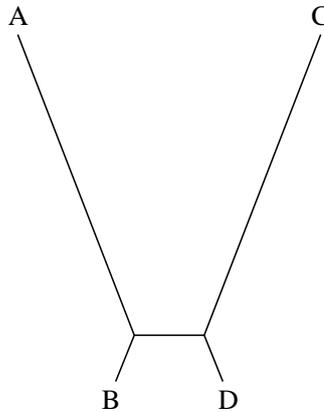
In all of the above discussion, we assumed that a multiple alignment of the sequences was given. If that is not the case, we have an even harder problem, the

**Generalized Tree Alignment Problem:** Given  $k$  sequences, find a tree with the sequences at the leaves and new reconstructed sequences at the internal nodes such that the length of the tree is minimal. Here the length of the tree is the sum of the edit distances of the sequences at the end of the edges.

This is a very hard problem. A dynamic programming algorithm that runs in exponential time exists [17]. Heuristic approximations also exist, see e.g. [22] and [23].

## 5.5 Long Branch Attraction

A standard objection to the parsimony criterion is that it is not *consistent*, i.e. even with infinite amount of data (infinitely long sequences), one will not necessarily obtain the correct tree. The standard example for this behavior is subsumed under the term *long branch attraction*. Assume the following tree with correct topology  $((A, B), (C, D))$ .



Since the branches pointing to  $A$  and  $C$  are very long, each of these two taxa looks essentially random when compared to the other three taxa. The only pair that does not look random is the pair  $(B, D)$ , and hence in a most parsimonious tree these two will be placed together giving the topology  $((B, D), (A, C))$ .

## **Part III**

# **Distance-based Methods**

## 6 Distance Based Trees

Distance based tree building methods rely on a distance measure between the taxa, resulting in a distance matrix. Distance measures usually take a multiple alignment of the sequences as input. After the distance measure is performed, sequence information is not used any more. This is in contrast to character based tree building methods which consider each column of a multiple sequence alignment as a character and which assess the nucleotides or amino acid residues at those sites (the character states) directly.

The idea when using distance based tree building methods is that knowledge of the “true evolutionary distances” between homologous sequences should enable us to reconstruct their evolutionary history.

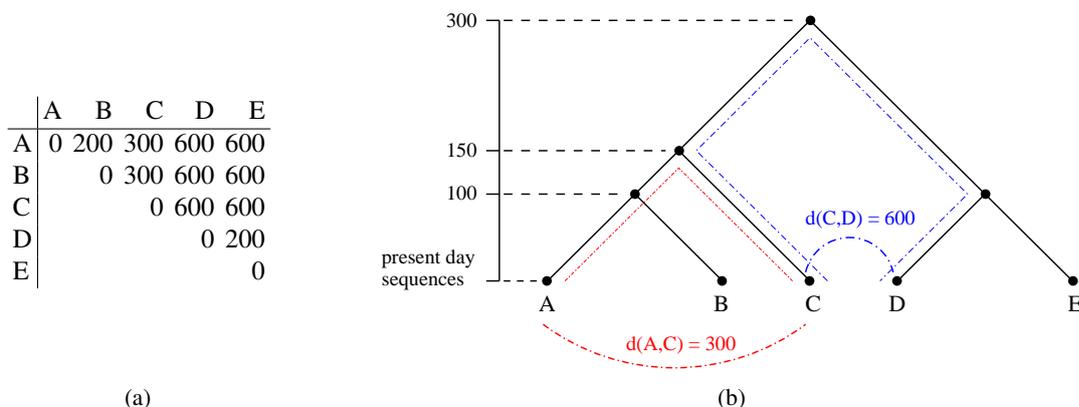
Suppose the evolutionary distances between members of a sequence set  $\{A, B, C, D, E\}$  were given by a distance matrix  $d^M$ , as shown in Figure 6.1(a).

Consider now the tree in Figure 6.1(b). A tree like this is called a *dendrogram* as the nodes are ranked on the basis of their relative distance to the root. The amount of evolution which has accumulated in A and C since divergence from their common ancestor is 150. In other words, the evolutionary distance from A (and C) to the common ancestor of A and C is 150. In general, the sum of edge weights along the path between two nodes corresponds to the evolutionary distance between the two nodes. Deriving distances between leaves is done by summing up edge weights along the path between the leaves. Distances derived in this way from a tree form the *path metric*  $d^T$  of the tree. For the data in Figure 6.1, we see that  $d^T = d^M$ . In general, it is not always possible to construct a “perfect” tree. In such a case, the aim is to find a tree whose path matrix  $d^T$  is as similar to the given matrix  $d^M$  as possible.

### 6.1 Basic Definitions

**Definition.** A *metric* on a set of objects  $O$  is given by an assignment of a real number  $d_{ij}$  (a distance) to each pair of objects  $(i, j)$ ,  $i, j \in O$ , if  $d_{ij}$  fulfills the following requirements:

- (i)  $d_{ij} > 0$  for  $i \neq j$
- (ii)  $d_{ij} = 0$  for  $i = j$

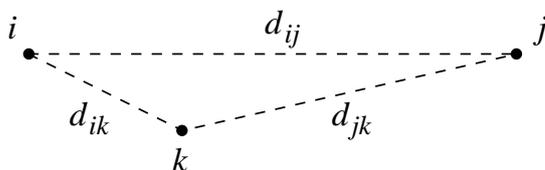


**Figure 6.1:** An example for distance based tree reconstruction methods. (a) The “true” evolutionary distances are given in a matrix  $d^M$ . For example, the value 300 in the first row of the above matrix shall be read as “the evolutionary distance between A and C is 300 units.” (b) Rooted phylogenetic tree with weighted edges (dendrogram). A root has been assigned assuming that the sequences have evolved from a common ancestor. The sum of edge weights along the path between two leaves is the evolutionary distance between the two leaves. These distances  $d^T$  correspond to the measured data, given in (a).

$$(iii) \quad d_{ij} = d_{ji} \quad \forall i, j \in O$$

$$(iv) \quad d_{ij} \leq d_{ik} + d_{kj} \quad \forall i, j, k \in O$$

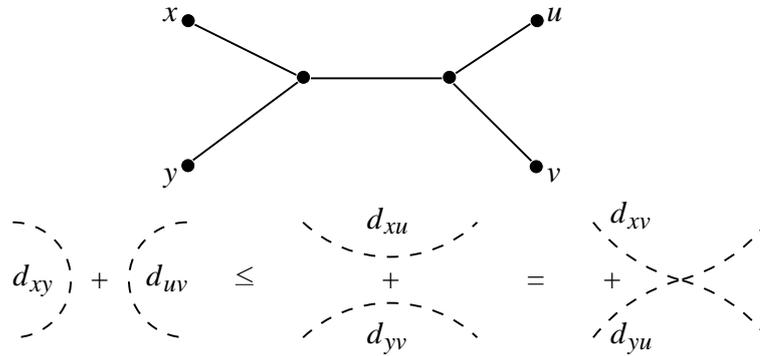
The latter requirement is called the *triangle inequality*:



**Definition.** Let  $d$  be a metric on a set of objects  $O$ , then  $d$  is an *additive metric* if it satisfies the inequality

$$d_{ij} + d_{kl} \leq \max(d_{ik} + d_{jl}, d_{il} + d_{jk}) \quad \forall i, j, k, l \in O.$$

An alternative and equivalent formulation is the *four point condition* according to Buneman [24]:



**Figure 6.2:** The four point condition. It implies that the path metric of a tree is an additive metric.

**Four point condition.** A metric  $d$  is an *additive metric* on  $O$ , if any four elements from  $O$  can be named  $x, y, u$  and  $v$  such that

$$d_{xy} + d_{uv} \leq d_{xu} + d_{yv} = d_{xv} + d_{yu}.$$

This is a stronger version of the triangle inequality. See Figure 6.2 for an illustration.

**Observation.** For an additive metric  $d^M$ , there exists a unique tree  $T$  such that  $d^T = d^M$ .

**Definition.** Let  $d$  be a metric on a set of objects  $O$ , then  $d$  is an *ultrametric* if it satisfies

$$d_{ij} \leq \max(d_{ik}, d_{jk}) \quad \forall i, j, k \in O.$$

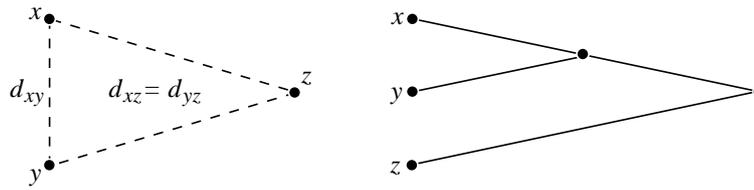
Again, there is an alternative formulation:

**Three point condition.** A metric  $d$  is an *ultrametric* on  $O$ , if any three elements from  $O$  can be named  $x, y, z$  such that

$$d_{xy} \leq d_{xz} = d_{yz}.$$

This is an even stronger version of the triangle inequality than the four point condition. If  $d$  is an ultrametric, it is an additive metric. See Figure 6.3 for an illustration.

**Observation.** For an ultrametric  $d^M$ , there exists a unique tree  $T$  that can be rooted in such a way that the distance from the root to any leaf is equal. Such a tree is called an *ultrametric tree*.



**Figure 6.3:** Three point condition. It implies that the path metric of a tree is an ultrametric.

## 6.2 Ultrametric Trees

The previous section dealt with the connection between matrices and trees. In the following, we discuss the reconstruction of ultrametric trees. Such a tree is given by the dendrogram in Figure 6.1(b). There is a clear interpretation inherent to ultrametric trees: Sequences have evolved from a common ancestor at constant rate (molecular clock hypothesis).

The path metric of an ultrametric tree is an ultrametric. Conversely, if distances  $d^M$  between a set of objects form an ultrametric, there is a unique ultrametric tree  $T$  corresponding to the distance measure, that is  $d^T = d^M$ . Note that  $T$  is not necessarily binary. A modified version of the three point condition ensures a unique binary tree:  $d_{xy} < d_{xz} = d_{yz}$ .

Given an ultrametric, the corresponding ultrametric tree can easily be reconstructed by one of the agglomerative clustering procedures described below.

Distance measures on real sequence data generally do not form an ultrametric. However, if the observed distances are close to an ultrametric, clustering procedures such as UPGMA (see Section 6.2.1) are the simplest and fastest way to reconstruct an ultrametric tree. While this is a very common approach, it is a heuristic (“algorithmic”) method which does not optimize a simple objective function. As mentioned above, being close to an ultrametric implies the existence of a molecular clock or a constant rate of evolution (which sometimes may hold for closely related sequences). Note that clustering procedures are sensitive to unequal evolutionary rates. If the assumption of rate constancy among lineages does not hold, UPGMA may give the wrong topology (for an example see Section 6.3.4 and Figure 6.4).

Other distance based methods like Neighbor Joining are more general, that is, they do not presume a molecular clock (see Section 6.3.4).

### 6.2.1 Agglomerative Clustering

Agglomerative clustering is conceptually simple and fast. Singleton clusters are successively merged into larger clusters to form a hierarchy:

Given a set of objects  $O$  with  $n$  elements and distances  $d_{ij}$ ,  $i, j \in O$ , initially each object is assigned a singleton cluster. Then the algorithm proceeds as follows:

While there is more than one cluster left, do:

1. Find the two clusters  $i$  and  $j$  with the smallest distance  $d_{ij}$ .
2. Create a new cluster  $u$  that joins clusters  $i$  and  $j$ .
3. Define the *height* (i.e. distance from leaves) of  $u$  to be  $l_{ij} := d_{ij}/2$
4. Compute the distance  $d_{ku}$  of  $u$  to any other cluster  $k \notin \{i, j\}$  in one of the ways described below.
5. Replace clusters  $i$  and  $j$  by the new cluster  $u$ .

Different clustering methods differ in how they define the distance  $d_{ku}$  between two clusters in step 4. For all variants presented here, a quadratic run time can be achieved.

*Single linkage clustering:*

$$d_{ku} := \min(d_{ki}, d_{kj}).$$

*Complete linkage clustering:*

$$d_{ku} := \max(d_{ki}, d_{kj}).$$

*UPGMA (unweighted pair group method using arithmetic averages):*

$$d_{ku} := \frac{1}{n_k n_u} \sum_{\substack{x \text{ object in } k, \\ y \text{ object in } u}} d_{xy}$$

where  $n_i$  is the number of objects in cluster  $i$ . The new distance  $d_{ku}$  is the arithmetic average of the original distances of all elements in  $k$  and all elements in  $u$ . A straight forward calculation would result in a cubic run time. However, the following formula yields the same distances and can be used to update them efficiently.

$$d_{ku} := \frac{n_i d_{ki} + n_j d_{kj}}{n_i + n_j} \tag{6.1}$$

UPGMA takes the single distances of the individual objects equally into account and is therefore called *unweighted*. The following method does not. It is thus called *weighted*.

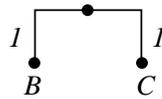
*WPGMA (weighted pair group method using arithmetic averages):*

$$d_{ku} := \frac{d_{ki} + d_{kj}}{2}.$$

**Example.** Given a set of objects  $O = \{A, B, C, D, E\}$  and an ultrametric distance matrix  $d^M$  on  $O$  with entries

	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<i>A</i>	0	8	8	12	8
<i>B</i>		0	2	12	4
<i>C</i>			0	12	4
<i>D</i>				0	12
<i>E</i>					0

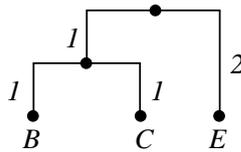
We want to reconstruct an ultrametric tree using UPGMA. As  $d_{BC} = 2$  is the smallest distance we join  $B$  and  $C$  into a new cluster  $(BC)$  with depth 1:



We update the distance matrix. E.g.  $d_{A(BC)} = (1 \cdot 8 + 1 \cdot 8)/(1 + 1) = 8$ , etc.

	<i>A</i>	<i>(BC)</i>	<i>D</i>	<i>E</i>
<i>A</i>	0	8	12	8
<i>(BC)</i>		0	12	4
<i>D</i>			0	12
<i>E</i>				0

We join  $(BC)$  and  $E$  with depth 2:

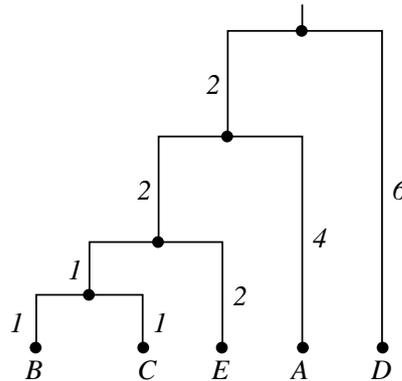


We obtain distances to  $((BC)E)$ , e.g.  $d_{A((BC)E)} = (2 \cdot 8 + 1 \cdot 8)/(2 + 1) = 8$ .

The modified distances are

	<i>A</i>	<i>((BC)E)</i>	<i>D</i>
<i>A</i>	0	8	12
<i>((BC)E)</i>		0	12
<i>D</i>			0

We join  $A$  and  $((BC)E)$  with depth 4, and finally  $(A((BC)E))$  and  $D$  are left to join:



Note that in this example single linkage, complete linkage and WPGMA yield the same unique tree. This is due to distances being ultrametric. If the data is not ultrametric, the results can be different.

UPGMA was originally developed for phenetics [25], i.e. for constructing phenograms reflecting phenotypic similarities rather than evolutionary distances. Given an approximately constant rate of evolution, that is if observed distances are close to an ultrametric, it is also suited for phylogeny reconstruction, and it is the most commonly used clustering method for this purpose.

Single linkage and complete linkage are somewhat extreme cases of clustering. While complete linkage clusters are usually very compact (each element in a cluster is connected to each other element), single linkage clusters may contain pairs of elements that by direct comparison are rather dissimilar when there is a “path” of connecting elements between them.

### 6.3 Additive Trees

We have seen that ultrametric trees are rooted trees and imply the existence of a molecular clock. But rates of evolution vary among species, among gene families, among sites in molecular sequences and generally in the course of sequence evolution. *Additive trees* do not presume a constant evolutionary rate nor do they make any assumption about the rooting and therefore reflect our ignorance as to where the common ancestor lies. Given an additive distance matrix there is exactly one tree topology that allows for realization of an additive tree. We will show how to reconstruct it in Section 6.3.1.

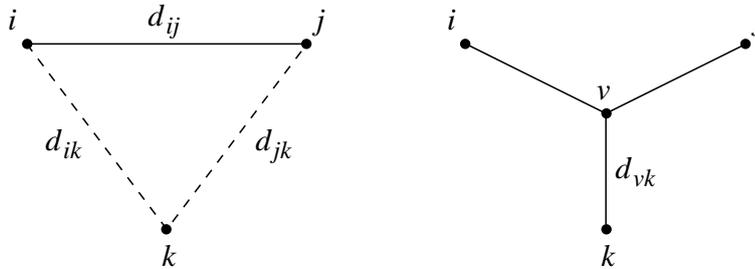
If one wants to construct a tree  $T$  from a distance matrix  $d^M$ , then the aim is that distances  $d^T$  are as similar as possible to the observed distances  $d^M$ . In Sections 6.3.2 and 6.3.3 we discuss two methods that optimize a simple objective function when observed distances  $d^M$  are not additive, the Fitch-Margoliash algorithm and the Minimum Evolution method, which aim at

reconstructing an additive tree  $T$  with distances  $d^T$  being as similar as possible to observed distances  $d^M$ . Finally in Section 6.3.4 we present the popular and heuristic method called Neighbor Joining.

### 6.3.1 Exact Reconstruction of Additive Trees

An additive metric can be represented as a unique additive tree which can be reconstructed in time complexity  $O(n^2)$  [26]. The algorithm successively inserts objects into intermediate trees until no objects are left to insert.

We use the following rationale: Given an intermediate tree  $T'$  containing leaf  $i$  and leaf  $j$ , we test if we can insert an edge connecting a new leaf  $k$  to the intermediate tree along the path between  $i$  and  $j$ . We denote the node connecting  $i$ ,  $j$  and  $k$  as  $v$  and the weight of the edge being inserted as  $d_{vk}$ .



We observe that

$$d_{ik} + d_{jk} = d_{iv} + d_{vk} + d_{jv} + d_{vk} = 2 \cdot d_{vk} + d_{ij}$$

and therefore the weight of the inserted edge would be

$$d_{vk} = \frac{1}{2}(d_{ik} + d_{jk} - d_{ij})$$

and respectively

$$\begin{aligned} d_{iv} &= d_{ik} - d_{vk} \\ d_{jv} &= d_{jk} - d_{vk}. \end{aligned}$$

The overall algorithm is then the following. Given a set of objects  $O$  and an additive metric  $d$  on  $O$ , one first picks two arbitrary objects  $i, j \in O$  and connects them by an edge with weight  $d_{ij}$ . This gives the first intermediate tree  $T'$ . Then, iteratively each object  $k \in O$  not yet in  $T'$  is connected to  $T'$  by an edge  $e_k$  by the following algorithm.

1. Pick a pair of leaves  $i, j \in T'$ .
2. Compute the weight of  $e_k$  by means of the above rationale.
3. If the insertion of  $e_k$  in  $T'$  implies that  $e_k$  has to be inserted inside an edge, split that edge, insert a node and attach  $e_k$  to that node; otherwise (if the insertion point is a node), replace  $j$  (or  $i$ ) by a leaf from the subtree along the edge at this node not leading towards  $i$  or  $j$  and continue with step 2.

To see that the algorithm runs in  $O(n^2)$  time, observe that there are  $n - 2$  iterations, each requiring a linear number of weight computations. The directed traversal of the tree in order to test if the branching point is a node or not can be done in time proportional to the number of nodes traversed with a few simple data structures.

Note that the algorithm can process non-additive metrics. But in this case, the resulting tree distances might be inconsistent with the given distances.

**Example:** Given a set of objects  $O = \{A, B, C, D, E\}$  and the following additive metric  $d$  on  $O$ :

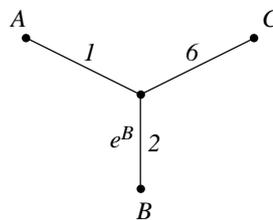
	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<i>A</i>	0	3	7	10	7
<i>B</i>		0	8	11	8
<i>C</i>			0	9	6
<i>D</i>				0	5
<i>E</i>					0

We first pick two arbitrary objects, say  $A$  and  $C$ , and connect them by an edge of weight  $d_{AC} = 7$  to set up the first intermediate tree  $T'$ :



We connect  $B$  by an edge  $e_B$  to  $T'$ . The weight of  $e_B$  is

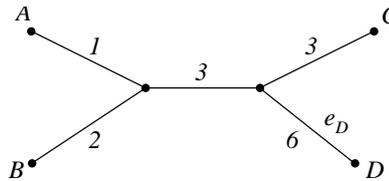
$$\frac{d_{AB} + d_{CB} - d_{AC}}{2} = \frac{3 + 8 - 7}{2} = 2.$$



We try to connect  $D$  by an edge  $e_D$  branching off the path between  $B$  and  $C$ . The weight of  $e_D$  would be

$$\frac{d_{BD} + d_{CD} - d_{BC}}{2} = \frac{11 + 9 - 8}{2} = 6$$

and inserting  $e_D$  on the edge branching off to  $C$  is therefore consistently possible:

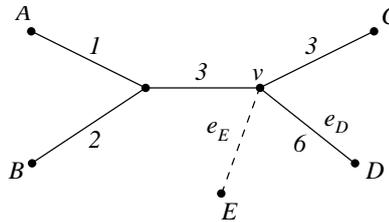


Finally we have to connect  $E$  by an edge  $e_E$ . We try to let  $e_E$  branch off the path between  $B$  and  $C$ . The weight of  $e_E$  would be

$$d_{vE} = \frac{d_{BE} + d_{CE} - d_{BC}}{2} = \frac{8 + 6 - 8}{2} = 3$$

and hence

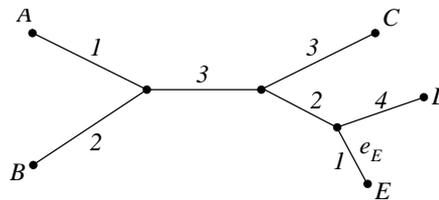
$$d_{Bv} = d_{BE} - d_{vE} = 8 - 3 = 5.$$



This implies that  $e_E$  has to be inserted at node  $v$ , and hence the procedure is repeated with  $C$  being replaced by  $D$ . Choosing the path between  $B$  and  $D$ , the weight of  $e_E$  is

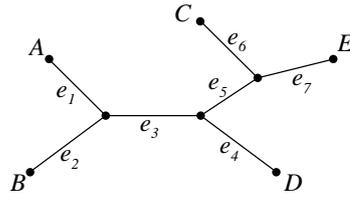
$$\frac{d_{BE} + d_{DE} - d_{BD}}{2} = \frac{8 + 5 - 11}{2} = 1$$

and as  $d_{DE} = 5$ ,  $e_E$  branches off  $e_D$ :



### 6.3.2 Least Squares (Fitch-Margoliash)

In practice, a distance measure  $d^M$  on a set of homologous sequences hardly will form an additive metric. If the given distances  $d^M$  are close to being additive, we want to find such a tree  $T$  with a path metric  $d^T$  that is as “similar” to  $d^M$  as possible. Fitch-Margoliash [27] and other methods are based on a family of objective functions called *least squares* to formally define “similarity”.



Consider the above tree  $T$  with its set of leaves  $O = \{A, B, C, D, E\}$  and weighted edges  $e_1, e_2, \dots, e_7$ . In general, for a tree with  $n$  leaves, let  $\vec{e}$  be the vector of the  $2n - 3$  edge weights  $w_j$ . In a binary tree, as shown below, we assign 1 to a pair of leaves  $(X, Y)$  and an edge  $e_j$  whenever  $e_j$  belongs to the simple path between  $X$  and  $Y$ , and 0 otherwise. Such a table, interpreted as a  $(\frac{n(n-1)}{2} \times 2n - 3)$  matrix  $M^T$ , is called the *path edge incidence matrix*. We can then see that

$$\vec{d}^T := M^T \vec{w}$$

is a vector holding the tree distances between the leaves of  $T$ .

$$\begin{pmatrix} d(A, B) \\ d(A, C) \\ d(A, D) \\ \vdots \\ d(C, E) \\ d(D, E) \end{pmatrix} = \begin{pmatrix} & e_1 & e_2 & e_3 & e_4 & e_5 & e_6 & e_7 \\ (A, B) & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ (A, C) & 1 & 0 & 1 & 0 & 1 & 1 & 0 \\ (A, D) & 1 & 0 & 1 & 1 & 0 & 0 & 0 \\ (A, E) & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ (B, C) & 0 & 1 & 1 & 0 & 1 & 1 & 0 \\ (B, D) & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ (B, E) & 0 & 1 & 1 & 0 & 1 & 0 & 1 \\ (C, D) & 0 & 0 & 0 & 1 & 1 & 1 & 0 \\ (C, E) & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ (D, E) & 0 & 0 & 0 & 1 & 1 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} w(e_1) \\ w(e_2) \\ \vdots \\ w(e_6) \\ w(e_7) \end{pmatrix}$$

Fitch and Margoliash [27] define the disagreement between a tree and the distance measure by

$$E := \|\vec{d}^T - \vec{d}^M\|^2 = \sum_i (d_i^T - d_i^M)^2.$$

One then wants to find a tree topology and edge lengths such that  $E$  is minimal. (To be precise, Fitch and Margoliash weight the difference of  $d^T$  and  $d^M$  by their measured distances, i.e.  $E := \sum_{i < j} (d_{ij}^T - d_{ij}^M)^2 / (d_{ij}^M)^2$ , because they want to minimize the square of the *relative error*, i.e. they assume that the uncertainty of the measurement is by the same *percentage* for all measurements.)

This can be solved by linear algebra: For a given tree topology (represented by a matrix  $M^T$ ), one can find the edge lengths  $\vec{w}$  that minimize

$$E = \|\vec{d}^T - \vec{d}^M\|^2 = \|M^T \vec{w} - \vec{d}^M\|^2,$$

e.g. by solving

$$M^T \vec{w} = \vec{d}^M$$

using singular value decomposition. A numerical solution is also possible and often faster.

This allows to optimize the branch lengths for a given topology. Still, one has to repeat this test for all (or many) topologies. The least squares tree is then the tree that minimizes  $E$ .

In practice this exact method takes very long. That is why Fitch and Margoliash suggest in their paper [27] a heuristic clustering algorithm to find tree topologies that have small (but not necessarily minimal) least squares error.

### 6.3.3 Minimum Evolution

Alternatively, Waterman *et al.* [26] have formulated a linear program to find a tree  $T$  with  $d^T$  similar to  $d^M$ . The objective function of the linear program is in spirit similar to the Steiner tree setup: The overall length of the tree is minimized. However, the least squares criterion is used only to fit the branch lengths, while to evaluate and compare trees, one uses the “tree length”

$$L := \sum_{i=1}^{2n-3} w_i$$

where the  $w_i$  are the  $2n - 3$  edge lengths, computed from the pairwise distances between the sequences as above. The tree minimizing  $L$  is called the *minimum evolution tree* (ME tree). There are two constraints to the linear program:

- (i) all the branch lengths are non-negative,

$$w_i \geq 0;$$

and

- (ii) the alignment distance between two sequences may underestimate the number of changes that have occurred between the two sequences in the course of evolution, but not overestimate. Therefore, for any pair of sequences, the tree distance is not allowed to be smaller than the measured distance:

$$d_{ij}^T \geq d_{ij}^M \quad \forall i, j.$$

Simulations have shown that Minimum Evolution consistently constructs better trees than Least Squares.

### Fast Minimum Evolution

In practice, Minimum Evolution can be implemented as follows [28]:

1. Apply Neighbor Joining (see Section 6.3.4) to get an initial tree. This takes  $O(n^3)$  time.
2. Refine the tree by branch swapping/nearest neighbor interchange, that is accepted whenever the amount of evolution decreases. This takes  $O(pn^3)$  time where  $p$  is the number of (tried) swaps.

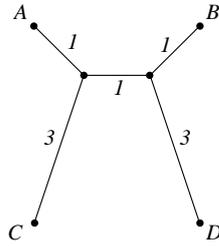
Desper and Gascuel [29] have suggested two speedups of the method:

1. A greedy addition algorithm that takes  $O(n^2)$  time.
2. A fast algorithm for nearest neighbor interchange, FASTNNI. This uses a clever pre-processing and takes only  $O(np)$  time where again  $p$  is the number of (tried) swaps.

### 6.3.4 Neighbor Joining

The *neighbor joining* (NJ) method is similar to cluster analysis in some ways. The individual taxa are iteratively grouped together, forming larger and larger clusters of taxa. In contrast to UPGMA, neighbor joining does not assume a molecular clock, but it assumes that observed distances are close to an additive metric. Given an additive metric, the neighbor joining method identifies the correct tree [30] and it also correctly reconstructs trees if additivity only holds approximately [31].

As neighbor relationships of nodes in a binary tree uniquely define the tree topology, successively identifying neighbors is a way to reconstruct the tree. In each iteration of the NJ algorithm, every pair of taxa is evaluated for being neighbors, and if so, they are grouped together to form a new taxon for the next iteration. Here, the notion of neighborhood is defined as follows:



**Figure 6.4:** Neighbors do not necessarily have the smallest distance.

**Definition.** Two taxa are *neighbors* in a tree if the path between them contains only one node.

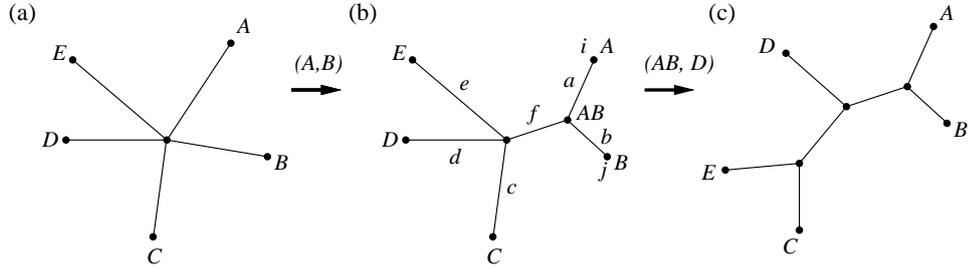
Note that neighbors do not need to have the smallest distance in the distance matrix. For an example, see the tree in Figure 6.4. The corresponding distance matrix  $d$  is:

	A	B	C	D
A	0	3	4	5
B		0	5	4
C			0	7
D				0

UPGMA would pick taxa  $A$  and  $B$  to cluster them together, since  $d_{AB}$  is the smallest distance. Actually  $A$  and  $B$  are not neighbors, but  $A$  and  $C$  are. This effect is due to the long edge with weight 3 (corresponding to a high rate at which mutations have accumulated) branching off to  $C$ . After presenting the central theorem and the algorithm, we show that Neighbor Joining correctly identifies the neighbor relationships for the tree in Figure 6.4.

The concept to identify neighbors is a variation of the Minimum Evolution principle (see also Section 6.3.3): A star tree is decomposed such that the tree length is minimized in each step. Consider the star tree with  $N$  leaves in Figure 6.5 (a). The star tree corresponds to the assumption that there is no clustering of taxa. In general there is a clustering of taxa and if so, the overall tree length (the sum of all branch lengths) of the true tree or the final NJ tree (see Figure 6.5 (c)) is smaller than the overall tree length of the star tree. Consider the tree in Figure 6.5 (b) with resolved neighbors  $A$  and  $B$ . It is clear that the tree length of this tree is smaller than that of the initial star tree. A general formula for the tree length  $S_{ij}$  of a tree like in Figure 6.5 (b) when considering taxa  $i$  and  $j$  as neighbors is

$$S_{ij} = \sum_{\substack{k=1 \\ k \neq i, j}}^N \frac{d_{ki} + d_{kj}}{2(N-2)} + \frac{d_{ij}}{2} + \sum_{\substack{k < l \\ k, l \neq i, j}}^N \frac{d_{kl}}{N-2}$$



**Figure 6.5:** Neighbor Joining successively decomposes a star tree by identifying neighbors. Pairs of neighbors are written in parenthesis.

where  $N$  is the number of taxa. Computation for the tree in Figure 6.5 (b) yields

$$S_{AB} = (3a + 3b + 6f + 2c + 2d + 2e) \cdot \frac{1}{6} + \frac{a + b}{2} + (2c + 2d + 2e) \cdot \frac{1}{3} = a + b + f + c + d + e$$

**Theorem:** Given an additive tree  $T$ .  $O$  is the set of leaves of  $T$ . Values of  $S_{ij}$  are computed by means of the path metric  $d^T$ . Then  $m, n \in O$  are neighbors in  $T$ , if  $S_{mn} \leq S_{ij}$  for all  $i, j \in O$ .

A simple proof makes use of the four point condition (see Section 6.1). It enables us to identify a pair of neighbors given additive distances between a set of taxa by computing  $S_{ij}$  for all pairs of taxa and choosing taxa  $i$  and  $j$  showing the smallest  $S_{ij}$  value. The identified neighbors are combined into one composite taxon and the procedure is repeated. We rewrite

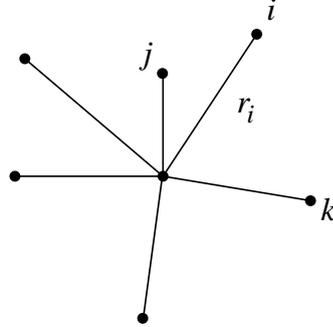
$$\begin{aligned} S_{ij} &= \frac{1}{2(N-2)} \left( 2 \sum_{\substack{k < l \\ k, l \neq i, j}}^N d_{kl} + \sum_{\substack{k=1 \\ k \neq i, j}}^N (d_{ki} + d_{kj}) \right) + \frac{d_{ij}}{2} \\ &= \frac{1}{2(N-2)} \left( 2 \sum_{k < l}^N d_{kl} - r_i - r_j \right) + \frac{d_{ij}}{2} \end{aligned}$$

with  $r_i := \sum_{k=1}^N d_{ik}$ . Since the sum  $\sum_{k < l}^N d_{kl}$  is the same for all pairs of  $i$  and  $j$ , we can replace  $S_{ij}$  by

$$M_{ij} := d_{ij} - \frac{r_i + r_j}{N - 2}$$

for the purpose of minimization:  $\operatorname{argmax}_{(i,j)}(S_{ij}) = \operatorname{argmax}_{(i,j)}(M_{ij})$ .

**Algorithm:** Given distances  $d_{ij}$  between members of a set  $O$  of  $N$  objects. Represent the objects as terminal nodes in a starlike tree:



1. For each terminal node  $i$  compute

$$r_i := \sum_{k=1}^N d_{ik}.$$

2. For all pairs of terminal nodes  $(i, j)$  compute

$$M_{ij} := d_{ij} - \frac{r_i + r_j}{N - 2}.$$

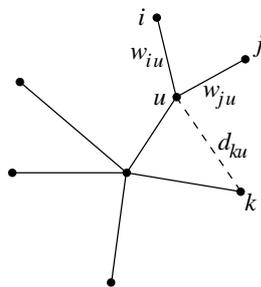
Let  $(i, j)$  be a pair with minimal value  $M_{ij}$  for  $i \neq j$ .

3. Join nodes  $i$  and  $j$  into a new terminal node  $u$ . The branch lengths from  $u$  to  $i$  and  $j$  are

$$w_{iu} = \frac{d_{ij}}{2} + \frac{r_i - r_j}{2N - 4} \quad \text{and} \quad w_{ju} = d_{ij} - w_{iu}.$$

4. Obtain the distances from  $u$  to another terminal node  $k$  by

$$d_{ku} = \frac{d_{ik} + d_{jk} - d_{ij}}{2}.$$



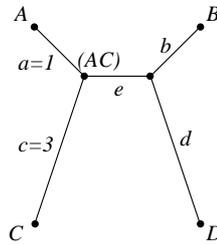
5. Delete  $i$  and  $j$  from the set of objects and replace them by  $u$ , thus reducing  $N$  by one. If there are more than two clusters left, continue with step 1.

Now we will give an example for a tree reconstruction given an (exact) additive metric by means of the NJ algorithm.

**Example.** The path metric  $d^T$  for the tree in Figure 6.4 is given by the following distances:

	$A$	$B$	$C$	$D$
$A$	0	3	4	5
$B$		0	5	4
$C$			0	7
$D$				0

In the first iteration  $A$ ,  $B$ ,  $C$  and  $D$  are the terminal nodes of the star tree and we compute  $r_A = r_B = 12$ ,  $r_C = r_D = 16$  and  $M_{AB} = d_{AB} - (r_A + r_B)/(N - 2) = 3 - 24/2 = -9$ ,  $M_{AC} = M_{BD} = 4 - 28/2 = -10$ ,  $M_{AD} = M_{BC} = 5 - 28/2 = -9$ ,  $M_{CD} = 7 - 32/2 = -9$ .  $M_{AC}$  and  $M_{BD}$  have the smallest value, that is, the NJ algorithm correctly identifies  $A$  and  $C$  as well as  $B$  and  $D$  as neighbors. We combine  $A$  and  $C$  into a composite taxon  $(AC)$ .



The edge lengths  $a$  and  $c$  are  $a = d_{AC}/2 + (r_A - r_C)/(2N - 4) = 2 + (-4/4) = 1$  and  $c = d_{AC} - a = 4 - 1 = 3$ . New distances of the composite taxon  $(AC)$  to  $B$  and  $D$  are  $d_{(AC)B} = (d_{AB} + d_{CB} - d_{AC})/2 = (3 + 5 - 4)/2 = 2$  and  $d_{(AC)D} = 4$ . We delete  $A$  and  $C$  from the set of objects and do the second iteration with the distance matrix

	$(AC)$	$B$	$D$
$(AC)$	0	2	4
$B$		0	4
$D$			0

There are three terminal nodes left and therefore we expect them all to be pairwise neighbors. Computations yield  $r_{(AC)} = 6$ ,  $r_B = 6$ ,  $r_D = 8$  and  $M_{(AC)B} = d_{(AC)B} - (r_{(AC)} + r_B)/(N - 2) = 2 - 12 = -10$ ,  $M_{(AC)D} = 4 - 14 = -10$ ,  $M_{(BD)} = -10$ . Grouping  $(AC)$  and  $B$  together into  $((AC)B)$ , we obtain  $e = d_{(AC)B}/2 + (r_{AC} - r_B)/2 = 1$  and  $b = 2 - 1 = 1$ . Now, there is only one distance left to compute:  $d = d_{(ACB)D} = (d_{(AC)D} + d_{BD} - d_{(AC)B})/2 = (4 + 4 - 2)/2 = 3$ . The NJ tree is the same as the true tree (Figure 6.4).

# 7 Phylogenetic Networks

## 7.1 Split Decomposition

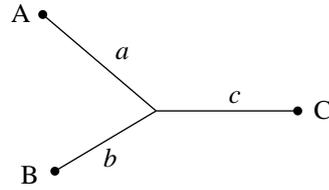
### 7.1.1 Introduction

The methods for reconstructing (ultrametric or additive) trees we have looked at before will always construct a *tree* even if the underlying distance data is not tree-like. In such a case a tree-like relationship will be suggested which is not present in the data. Unfortunately in most cases the methods do not even tell how little the data reflects the resulting tree. And even if they do (as in the least-squares method), only a single quantity is returned without any information as to specific regions where the dissimilarity occurs or what alternative trees might be.

A method called *split decomposition* developed by Bandelt and Dress [32, 33] allows both to quantify the tree-likeness of given distance data and present alternative relationships. A measured dissimilarity matrix  $d^M$  is decomposed into a number of *splits* (binary partitions of the set of taxa) weighted by *isolation indices* (indicating the strength of the split), plus a residual *noise* term. This is motivated by the assumption that measured distance data may underly a systematic error. Assume that  $d^M = d^T + d^E$  where  $d^T$  is the true tree-like relationship and  $d^E$  is an error term. If  $d^E$  itself represents a tree (different from the true one), then the two sets of splits will overlay. Ab initio it will not be possible to distinguish which splits are the true ones, and which result from the error term, but if the true splits are stronger than the error splits, one can assume that those splits with the larger isolation index belong to  $d^T$  rather than to  $d^E$ .

### 7.1.2 Basic Idea

First remember the case of three taxa A, B, C and three distances  $d_{AB}$ ,  $d_{AC}$ ,  $d_{BC}$ . As long as the triangle inequality ( $d_{AC} \leq d_{AB} + d_{BC}$ ) holds, it is always possible to compute a unique tree. We have three conditions and three unknowns (the lengths of the three terminal edges)  $a$ ,  $b$ ,  $c$ :

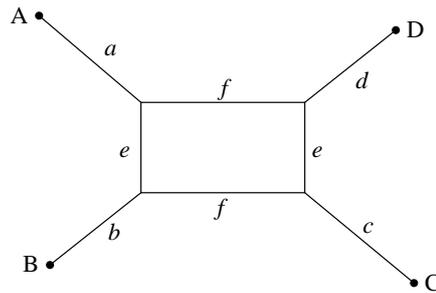


$$a = \frac{1}{2}(d_{AB} + d_{AC} - d_{BC}),$$

$$b = \frac{1}{2}(d_{AB} + d_{BC} - d_{AC}),$$

$$c = \frac{1}{2}(d_{AC} + d_{BC} - d_{AB}).$$

In the case of four taxa, a tree is no longer uniquely defined: we have six conditions, but a tree with four leaves has only five edges. However, the following diagram, which shows a generalized “tree” whose internal edges are replaced by a rectangle, has six unknown edge lengths  $a, b, c, d, e, f$ :



$$a = \frac{1}{2}(d_{AB} + d_{AD} - d_{BD}),$$

$$b = \frac{1}{2}(d_{AB} + d_{BC} - d_{AC}),$$

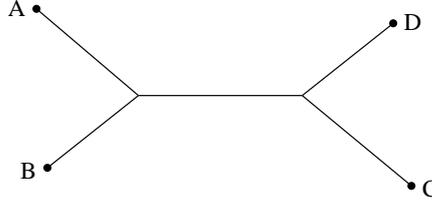
$$c = \frac{1}{2}(d_{BC} + d_{CD} - d_{BD}),$$

$$d = \frac{1}{2}(d_{AD} + d_{CD} - d_{AC}),$$

$$e = \frac{1}{2}(d_{AC} + d_{BD} - d_{AD} - d_{BC}),$$

$$f = \frac{1}{2}(d_{AC} + d_{BD} - d_{AB} - d_{DC}).$$

Now assume that the phylogenetic relationship is such that the true split separates the pair (A,B) from the pair (C,D).



Then due to the four-point condition the true tree distances  $d^T$  are related as follows:

$$d_{AB}^T + d_{CD}^T \leq d_{AC}^T + d_{BD}^T = d_{AD}^T + d_{BC}^T.$$

As discussed above, the *measured* (evolutionary) distances  $d^M$  will usually not fulfill this relationship, normally not even the two orderings

$$d_{AB}^M + d_{CD}^M \leq d_{AC}^M + d_{BD}^M \quad \text{and} \quad d_{AB}^M + d_{CD}^M \leq d_{AD}^M + d_{BC}^M.$$

However, one could hope that at least  $d_{AB}^M + d_{CD}^M$  is not the largest of the three sums.

### 7.1.3 Definitions

Based on this idea, one can define a *split with respect to  $d^M$*  or a  $d^M$ -*split* (or  $d$ -split when it is clear from the context that we mean  $d^M$ ) between a group  $J$  and a group  $K$  of taxa as follows. For any two taxa  $i, j$  from  $J$  and  $k, l$  from  $K$ , the sum  $d_{ij} + d_{kl}$  must not be the largest among the three possible sums, i.e.

$$d_{ij} + d_{kl} \leq \max\{d_{ik} + d_{jl}, d_{il} + d_{jk}\}.$$

It can be shown that for  $n$  taxa there can be at most  $\binom{n}{2}$  such splits [32, Theorem 3 and Corollary 4]. This number, however, is considerably larger than  $2n - 3$ , the maximal number of splits (edges) in an additive tree. Hence, splits define a relationship more general than a tree, though the number of  $d$ -splits observed in real data is usually only about  $2n$ .

Each  $d$ -split receives a positive weight. The *isolation index*  $\alpha_{J,K}$  for  $J, K$  is defined as

$$\alpha_{J,K} = \frac{1}{2} \min_{\substack{i,j \in J \\ k,l \in K}} (\max\{d_{ij} + d_{kl}, d_{ik} + d_{jl}, d_{il} + d_{jk}\} - d_{ij} - d_{kl}).$$

Note that all partitions that do not qualify as  $d$ -splits have isolation index 0. Moreover, the isolation index of a split in an additive tree is the length of the edge that defines the split.

The *split metric*  $\delta_{J,K}$  defined by a split  $J, K$  assigns distance 0 to two taxa if both are in  $J$  or both are in  $K$ , and 1 otherwise. Moreover, define  $d^1$  to be the sum of all split metrics weighted by their isolation indices,

$$d^1 = \sum_{d\text{-splits } J,K} \alpha_{J,K} \delta_{J,K}.$$

Then it can be shown [32] that  $d$  can be approximated by  $d = d^0 + d^1$  where  $d^0$  is a metric that does not contain any further splits. The proportion of  $d$  that can be split is called the *splittable percentage*  $\rho$  where

$$\rho := \left( \frac{\sum_{\text{taxa } i,j} d_{ij}^1}{\sum_{\text{taxa } i,j} d_{ij}} \right) \cdot 100\%.$$

#### 7.1.4 Computation of the $d$ -Splits

The set of  $d$ -splits of a distance matrix of  $n$  taxa  $1, 2, \dots, n$  can be computed by the following recursive algorithm:

1. Start with taxa 1, 2, 3, 4. This case can easily be solved exhaustively.
2. For  $i = 5, \dots, n$  do the following:
  - (i) Assume the  $d$ -splits restricted to the subset  $\{1, \dots, i - 1\}$  are given. For each  $d$ -split  $J, K$  of this subset, test if  $J \cup \{i\}, K$  or  $J, K \cup \{i\}$  qualifies as a  $d$ -split.
  - (ii) Test if  $\{1, \dots, i - 1\}, \{i\}$  is a  $d$ -split.

The algorithm can easily be implemented such that it has an  $O(n^6)$  time complexity.

**Graph drawing.** For a tree this is trivial: splits correspond to edges. In general, the resulting graphs are subgraphs of the  $\binom{n}{2}$ -dimensional hypercube, but usually they are rather simple, often planar. A split corresponds to several parallel edges, the length of which is proportional to the isolation index of the split.

The splits diagram can be constructed incrementally. The outcome is not unique; it depends on the order in which the splits are processed.

**Tree selection.** In addition, one may choose a set of splits that define a tree by greedily selecting those splits with maximal isolation index if they are compatible (every pair has parts with empty intersection) with previously selected ones.

**Program.** The splits method is implemented in a program called *Splitstree* that can be downloaded from <http://www.splitstree.org>. It finds the splits, draws a picture of the splits diagram, and calculates the splittable percentage. An example from the SplitsTree home page is given in Figure 7.1.

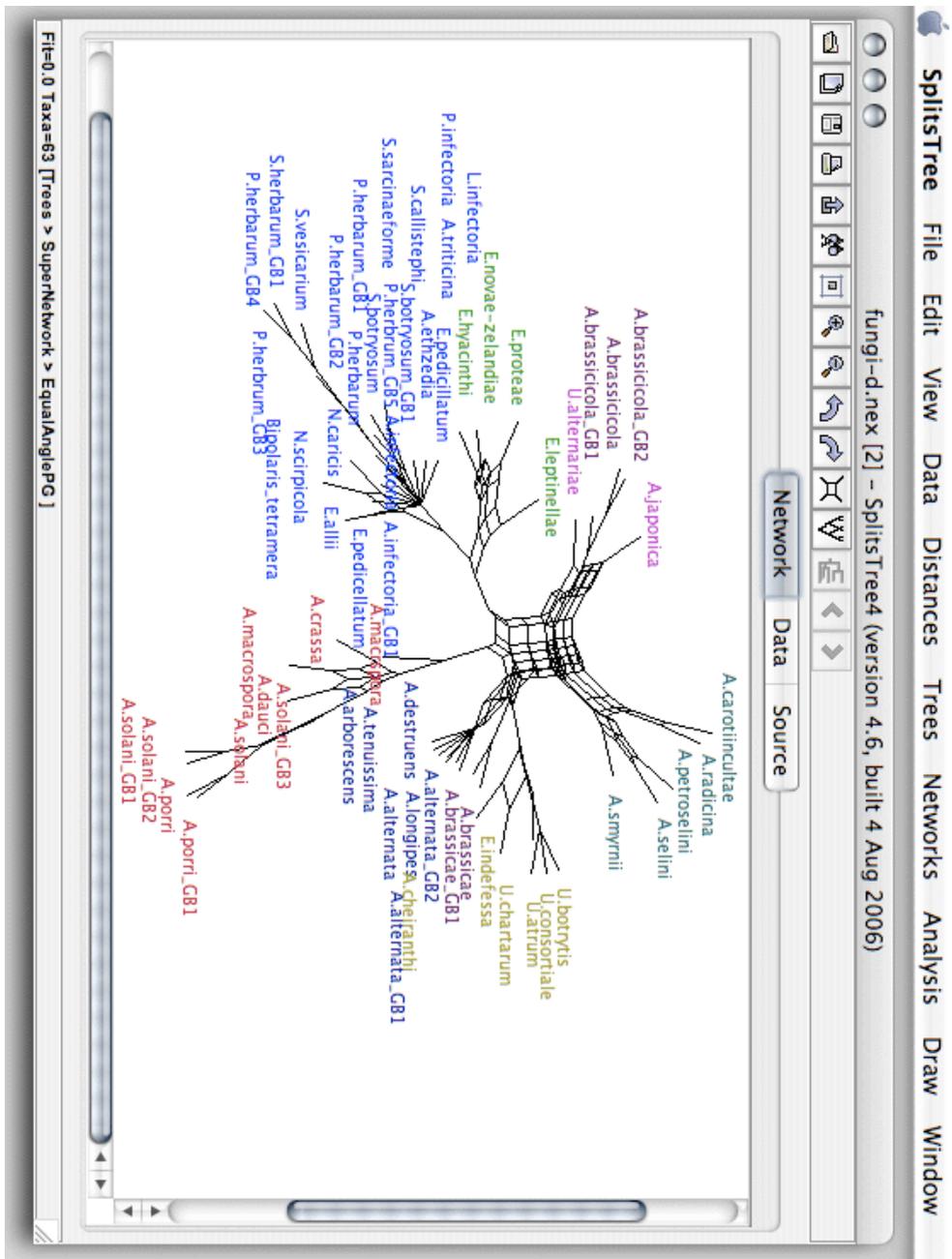


Figure 7.1: SplitsTree4 screen shot.

**Case studies.** Bandelt and Dress [33] presented a number of interesting case studies that illustrate the splits method.

## 7.2 NeighborNet

Recently, a few other methods for reconstructing phylogenetic networks have been developed. One is NeighborNet by Bryant and Moulton [34] which is a combination of Neighbor Joining and SplitsTree. The network is constructed incrementally, similar to Neighbor Joining. The result is a network as in SplitsTree, but with a better resolution. The running time is  $O(n^3)$ .

**Part IV**

**Likelihood Methods**

# 8 Modeling Sequence Evolution

## 8.1 Basics on Probability

### 8.1.1 Events and Probabilities

Throwing a die may be seen as an experiment: We do not know in advance which number will turn up. All we know is that a number between 1 and 6 will turn up.

**Definition.** The set of all possible outcomes of an experiment is called *sample space*  $\Omega$ . A subset  $A \subseteq \Omega$  is called an *event*. Elements  $\omega \in \Omega$  are called *elementary events*. The event  $\Omega$  is called the *certain event* and the event  $\emptyset$  is called the *null event*. The *complement* of an event  $A$  is denoted by  $A^C := \Omega \setminus A$ . Events  $A$  and  $B$  are called *disjoint* if  $A \cap B = \emptyset$ .

**Example.** A die is thrown. The sample space is  $\Omega = \{1, 2, 3, 4, 5, 6\}$ . The event  $A = \{2, 4, 6\}$  means that the outcome of the experiment is an even number.

We want to assign real numbers representing probabilities to subsets of  $\Omega$ , that is to events. If  $\Omega$  is infinite, e.g.  $\Omega = \mathbb{R}$ , it is not possible to reasonably assign probabilities to all members of the *power set* denoted by  $\{0, 1\}^\Omega$  which contains all subsets of  $\Omega$ . Therefore, we introduce the collection  $\mathcal{F}$  of subsets of  $\Omega$  containing all events of interest.

**Definition.** A collection  $\mathcal{F}$  of subsets of  $\Omega$  is called a  *$\sigma$ -field* if it satisfies the following conditions:

- (i)  $\emptyset \in \mathcal{F}$ ,
- (ii)  $A_1, A_2, \dots \in \mathcal{F} \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$ ,
- (iii)  $A \in \mathcal{F} \Rightarrow A^C \in \mathcal{F}$ .

**Example.** A die is thrown. The sample space is  $\Omega = \{1, 2, 3, 4, 5, 6\}$ .  $\mathcal{F} = \{\emptyset, \{2, 4, 6\}, \{1, 3, 5\}, \Omega\}$  is a  $\sigma$ -field.

We proceed by defining the properties of a function  $Pr$  assigning probabilities to events.

**Definition.** A probability measure  $Pr$  on  $(\Omega, \mathcal{F})$  is a function  $Pr : \mathcal{F} \rightarrow [0, 1]$  satisfying

- (i)  $Pr(\emptyset) = 0, \quad Pr(\Omega) = 1;$
- (ii) if  $A_1, A_2, \dots$  is a collection of *disjoint* members of  $\mathcal{F}$  then

$$Pr\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} Pr(A_i).$$

The triple  $(\Omega, \mathcal{F}, Pr)$  is called a *probability space*,  $Pr(A)$  is called the *probability* of event  $A$ .

**Example.** A die is thrown. The sample space is  $\Omega = \{1, 2, 3, 4, 5, 6\}$  and we can take  $\mathcal{F} = \{0, 1\}^{\Omega}$ . The elementary event  $\omega_i$  denotes the event that number  $i$  turns up. Suppose the die is fair, that is, each elementary event  $\omega_i$  has the same chance to occur. As  $Pr(\Omega) = 1$  and  $\bigcup_{i=1}^6 \omega_i = \Omega$  where  $\omega_i \cap \omega_j = \emptyset$  for all  $i, j \in \Omega$  with  $i \neq j$ , we see that  $Pr(\omega_i) = p = 1/6$ . The probability for the event  $A = \{\omega_1, \omega_2\}$  that ‘1’ or ‘2’ turns up is  $Pr(A) = 2p = 1/3$ .

### 8.1.2 Conditional Probability

Suppose, we have some prior knowledge of the outcome of an experiment.

**Definition.** The *conditional probability* that an event  $A$  occurs given another event  $B$  is

$$Pr(A|B) = \frac{Pr(A \cap B)}{Pr(B)}, \quad \text{if } Pr(B) > 0.$$

**Example.** A fair die is thrown. What is the probability for the event  $A$  that the number turning up is bigger than 1 given we know that the number is even? The probability for the event  $B$  that the number is even is  $Pr(B) = 1/2$ . We see that  $Pr(A \cap B) = Pr(\{2, 3, 4, 5, 6\} \cap \{2, 4, 6\}) = Pr(\{2, 4, 6\}) = 1/2$ . Therefore we get  $Pr(A|B) = (1/2)/(1/2) = 1$ .

**Lemma.** If  $A$  and  $B$  are events and  $Pr(B) > 0$  and  $Pr(B^C) > 0$  then

$$Pr(A) = Pr(A|B) Pr(B) + Pr(A|B^C) Pr(B^C).$$

**Proof.**  $A = (A \cap B) \cup (A \cap B^C)$ . As  $(A \cap B) \cap (A \cap B^C) = \emptyset$  we get  $Pr(A) = Pr(A \cap B) + Pr(A \cap B^C) = Pr(A|B)Pr(B) + Pr(A|B^C)Pr(B^C)$ .  $\square$

**Example.** A fair die is thrown. Let  $A$  be the event that the number turning up is even and  $B$  the event that the number is greater than two. Then  $Pr(B) = 2/3$  and  $Pr(B^C) = 1/3$ . As expected, we get  $Pr(A) = Pr(A|B)Pr(B) + Pr(A|B^C)Pr(B^C) = 1/2 \cdot 2/3 + 1/2 \cdot 1/3 = 1/2$ .

### 8.1.3 Bayes's Formula

From the definition of conditional probability we obtain

$$Pr(A \cap B) = Pr(A|B) Pr(B) = Pr(B \cap A) = Pr(B|A) Pr(A).$$

Solving for  $Pr(A|B)$  we obtain *Bayes's formula*

$$Pr(A|B) = \frac{Pr(B|A) Pr(A)}{Pr(B)}$$

which often turns out to be useful for the computation of conditional probabilities.

**Example.** A fair die is thrown. Let  $A$  be the event that the number turning up is greater than one and  $B$  the event that the number is even. Then  $Pr(A) = 5/6$ ,  $Pr(B) = 1/2$  and  $Pr(B|A) = 3/5$ . Using Bayes's formula, we get  $Pr(A|B) = (3/5 \cdot 5/6)/(1/2) = 1$ .

### 8.1.4 Independence

**Definition.** Two events  $A$  and  $B$  are *independent* if

$$Pr(A \cap B) = Pr(A)Pr(B).$$

**Example.** A fair die is thrown. Suppose event  $A$  is that the number turning up is bigger than 2 and event  $B$  is that the number turning up is even.  $A \cap B = \{4, 6\}$  and therefore  $Pr(A \cap B) = 1/3$ . We see that  $Pr(A)Pr(B) = 2/3 \cdot 1/2 = \frac{1}{3} = Pr(A \cap B)$ . Events  $A$  and  $B$  are independent.

**Example.** Two fair dice are rolled. The sample space is  $\Omega = \{(a, b) : a = 1, \dots, 6; b = 1, \dots, 6\}$  and  $|\Omega| = 36$ . Let events  $A$  and  $B$  be that a ‘1’ turns up on the first and the second die, respectively. The probability for the event  $(1, 1)$  that two times ‘1’ turns up is  $Pr(A \cap B) = 1/6 \cdot 1/6 = 1/36$ .  $A$  and  $B$  are independent.

### 8.1.5 Random Variables

Often we are not directly interested in the outcome of an experiment but in a function on the outcome, e.g. in the sum of numbers when rolling two dice.

**Definition.** A *random variable* is a function  $X : \Omega \rightarrow \mathcal{X}$  where  $\mathcal{X}$  may be any set.

**Example.** Two fair dice are rolled. The sample space is  $\Omega = \{(a, b) : a = 1, \dots, 6; b = 1, \dots, 6\}$ . We define the random variable  $X$  by  $X(\omega) := a + b$  for  $\omega = (a, b) \in \Omega$ . Therefore  $\mathcal{X} = \{2, 3, \dots, 12\}$ . The probability that  $X$  takes value ‘3’ is  $Pr(X = 3) = Pr(\{(1, 2)\}) + Pr(\{(2, 1)\}) = 1/18$ .

## 8.2 Markov Chains

Assuming that sites in DNA and amino acid sequences evolve independently from each other, evolution of molecular sequences is commonly modeled by means of a Markov chain at each site.

### 8.2.1 Time Discrete Markov Chains

**Definition.** A *time discrete Markov chain* is a sequence of random variables  $X_n$ ,  $n \in \mathbb{N}_0$ , taking values of a finite *set of states*  $\mathcal{A}$  where

- (i)  $X_0$  is sampled according to an *initial distribution*  $\pi^{(0)}$ :

$$\pi_i^{(0)} = Pr(X_0 = x_i), \quad x_i \in \mathcal{A};$$

- (ii) the *Markov property* has to be satisfied:

$$Pr(X_n = x_n | X_{n-1} = x_{n-1}, \dots, X_0 = x_0) = Pr(X_n = x_n | X_{n-1} = x_{n-1}),$$

for all  $n \in \mathbb{N}$  and all states  $x_0, \dots, x_n \in \mathcal{A}$ .

When modeling sequence evolution,  $\mathcal{A}$  will usually be the set of amino acid residues  $\mathcal{A} = \{A, C, \dots, Y\}$  or the set of nucleotides  $\mathcal{A} = \{A, C, G, T\}$ , respectively. Thus we can think of the Markov Chain as describing the behavior of a site in a molecular sequence in time. The Markov property means that the conditional probability for the observation of a state at a given time point only depends on the previous time point. The conditional probability for a state to reach another state is called *transition probability*.

We assume that the Markov chain is *time homogeneous*, that is for each pair of states  $x_i, x_j \in \mathcal{A}$  the transition probability  $Pr(X_n = x_j | X_{n-1} = x_i)$  is the same for all  $n \in \mathbb{N}$ .

Transition probabilities are held in the *transition probability matrix*  $P$  with entries

$$P_{ij} = Pr(X_n = x_j | X_{n-1} = x_i), \quad x_i, x_j \in \mathcal{A}.$$

$P$  is a stochastic matrix, i.e. each entry is nonnegative,  $P_{ij} \geq 0$  for all  $x_i, x_j \in \mathcal{A}$ , and each row sums up to 1,  $\sum_j P_{ij} = 1$ . The pair  $(\pi^{(0)}, P)$  specifies a unique homogeneous Markov chain.

We want to compute transition probabilities for  $k$  steps of the Markov chain. Let's start with two steps and  $\mathcal{A} = \{A, C, G, T\}$ . We obtain

$$\begin{aligned} Pr(X_{n+1} = T | X_{n-1} = A) &= Pr(X_n = A | X_{n-1} = A) Pr(X_{n+1} = T | X_n = A) \\ &\quad + Pr(X_n = C | X_{n-1} = A) Pr(X_{n+1} = T | X_n = C) \\ &\quad + Pr(X_n = G | X_{n-1} = A) Pr(X_{n+1} = T | X_n = G) \\ &\quad + Pr(X_n = T | X_{n-1} = A) Pr(X_{n+1} = T | X_n = T) \\ &= \sum_{k \in \mathcal{A}} P_{1k} P_{k4}. \end{aligned}$$

We denote  $P^{(1)} = P^1 = P$  as the 1-step transition matrix. We see that the 2-step transition matrix is  $P^{(2)} = P^2 = P \cdot P$ , and generally the  $k$ -step transition matrix is  $P^{(k)} = P^k = P P^{k-1}$ .

**Definition.** A Markov chain is *irreducible* if for any two states  $x_i, x_j \in \mathcal{A}$  there exists a  $k \in \mathbb{N}$  such that  $P_{ij}^{(k)} > 0$ .

That is, a Markov chain is irreducible if it is possible for the chain to reach each state from each state. This generally holds for the Markov chains in our applications.

The distribution of the states after  $k$  time steps,  $\pi^{(k)}$ , is a row vector that is obtained from the initial distribution  $\pi^{(0)}$  by

$$\begin{aligned}\pi_j^{(k)} &= \sum_{x_j \in \mathcal{A}} \pi_j^{(0)} \cdot P_{ij}^{(k)}, \\ \pi^{(k)} &= \pi^{(0)} \cdot P^{(k)}.\end{aligned}$$

**Definition.** A distribution  $\pi$  is a *stationary distribution* on  $\mathcal{A}$  if

$$\begin{aligned}\pi_j &= \sum_{x_i \in \mathcal{A}} \pi_i P_{ij} \quad \forall x_j \in \mathcal{A}, \\ \pi &= \pi P.\end{aligned}$$

If  $\pi^{(0)} = \pi$ , every  $X_n$  is distributed as  $\pi$ . Furthermore each irreducible homogeneous Markov chain converges against its stationary distribution:

**Theorem.** Given an irreducible time homogeneous Markov Chain  $(\pi^{(0)}, P)$ , there exists exactly one stationary distribution  $\pi$  and

$$\lim_{k \rightarrow \infty} P_{ij}^{(k)} = \pi_j \quad \forall x_i \in \mathcal{A}.$$

## 8.2.2 Time Continuous Markov Chains

In this subsection we will transfer the notions from time discrete Markov chains to time continuous Markov chains.

**Definition.** A *time continuous Markov chain* is a sequence of random variables  $X_t$ ,  $t \in \mathbb{R}_0^+$  taking values of a finite set of states  $\mathcal{A}$ .  $X_{t_0}$  is distributed as  $\pi^{(0)}$  and the Markov property holds:

$$Pr(X_{t_n} = x_n | X_{t_{n-1}} = x_{n-1}, \dots, X_{t_0} = x_0) = Pr(X_{t_n} = x_n | X_{t_{n-1}} = x_{n-1})$$

for all  $n \in \mathbb{N}$ , time points  $t_0 < t_1 < \dots < t_n$  and all states  $x_0, x_1, \dots, x_n \in \mathcal{A}$ .

The Markov chain is *time homogeneous* if there exists a *transition probability matrix*  $P(t)$  such that

$$Pr(X_{s+t} = x_j | X_s = x_i) = P_{ij}(t), \quad \forall s, t \geq 0, \quad \forall x_i, x_j \in \mathcal{A}.$$

The transition probability matrix  $P(t)$  is a stochastic matrix and has the following properties:

- $P(0) = \mathbb{1}$  (where  $\mathbb{1}$  is the identity matrix),
- $P_{ij}(t) \geq 0$  and  $\sum_{x_j \in \mathcal{A}} P_{ij}(t) = 1$  ( $P$  is a stochastic matrix),
- $P(s+t) = P(s)P(t)$  for  $s, t \geq 0$  (*Chapman-Kolmogorov equation*).

A time continuous Markov chain is *irreducible* if for any period  $t > 0$  each state can reach each state:  $P_{ij}(t) > 0$  for all  $x_i, x_j \in \mathcal{A}$ . In that case there exists a unique *stationary distribution*  $\pi$  which is the solution of  $\pi P(t) = \pi$ .

### 8.2.3 The Rate Matrix

We assume that the transition probability matrix  $P(t)$  of a time continuous Markov chain is continuous and differentiable at any  $t > 0$ . I.e. the limit

$$\lim_{t \searrow 0} \frac{P(t) - \mathbb{1}}{t} = Q$$

exists.  $Q$  is known as the *rate matrix* or, alternatively, the *generator* of the Markov chain. For very small time periods  $h > 0$ , transition probabilities are approximated by

$$\begin{aligned} P(h) &\approx \mathbb{1} + hQ \\ P_{ij}(h) &\approx Q_{ij} \cdot h, \quad i \neq j. \end{aligned}$$

From the last equation we see that the entries of  $Q$  may be interpreted as substitution rate per site per year.

From the Chapman-Kolmogorov equation we get the forward and backward equation

$$\begin{aligned} \frac{d}{dt}P(t) &= \lim_{h \searrow 0} \frac{P(t+h) - P(t)}{h} \\ &= \lim_{h \searrow 0} \frac{P(t)P(h) - P(t)\mathbb{1}}{h} \\ &= P(t) \lim_{h \searrow 0} \frac{P(h) - P(0)}{h} \\ \frac{d}{dt}P(t) &= P(t)Q = QP(t). \end{aligned}$$

This differential equation can be solved under the initial condition  $P(0) = \mathbb{1}$  and yields

$$P(t) = \exp(tQ) = \sum_{k=0}^{\infty} \frac{Q^k t^k}{k!}.$$

Thus, transition probabilities for any time  $t > 0$  may be computed from the matrix  $Q$ .  $Q$  provides an infinitesimal description of the process. As  $\sum_j P_{ij}(t) = 1$  we have

$$\sum_j Q_{ij} = 0$$

(the rows of the rate matrix sum to 0) and therefore  $Q_{ij} \geq 0$  for  $i \neq j$  and  $Q_{ii} \leq 0$ .

We may denote the time continuous Markov chain by  $(\pi^{(0)}, Q)$  (if  $Q$  exists and has the above properties). Then  $\pi$  is a stationary distribution if

$$\begin{aligned} \sum_i \pi_i Q_{ij} &= 0 \quad \forall j, \\ \pi Q &= \vec{0}. \end{aligned}$$

### 8.2.4 Definition of an Evolutionary Markov Process (EMP)

So far we have collected basic notions on time discrete and time continuous Markov chains. We additionally make the assumption that  $X_t$  is *reversible* and calibrate the rate matrix  $Q$  to a time unit. This enables us to define an *evolutionary Markov process (EMP)* according to Müller and Vingron [35]. The definition of an EMP summarizes the requirements on the Markov chain such that it is suited to describe the substitution process at a site of a molecular sequence. We start by defining time units.

**Definition.** One *PEM* (percent of expected mutations) is the time in which one substitution event (mutation) per 100 sites is expected in. One *PAM* (percent of accepted mutations) is the time for which the average number of substitutions per 100 sites is one.

- Given a rate matrix  $Q$ , one expects per time unit

$$E := \sum_i \pi_i \sum_{j \neq i} Q_{ij} = \underbrace{\sum_i \pi_i Q_{ij}}_{=0} - \sum_i \pi_i Q_{ii} = - \sum_i \pi_i Q_{ii}$$

substitutions per site. We calibrate  $Q$  by multiplying it with a constant. If  $E = 1/100$ ,  $Q$  is calibrated to 1 PEM.

- Given a 1-step probability transition matrix  $P^{(1)}$

$$E' := \sum_i \pi_i \sum_{j \neq i} P_{ij}^{(1)} = 1 - \sum_i \pi_i P_{ii}^{(1)}$$

is the average number of sites being in another state after one time unit. If  $E' = 1/100$ ,  $P^{(1)}$  is calibrated to 1 PAM.

In contrast to PAM units, PEM units take back-mutations into account. Therefore, 1 PEM is a slightly shorter time unit than 1 PAM.

We observe pairs of homologous sequences having evolved from a common ancestor. Yet we do not have any information about ancestral sequences. Therefore, we assume that evolution from ancestors to descendants can be modeled by the same process as its reverse. Thus the divergence of two homologous present-day sequences is explained by one process. This property of Markov chains is called *reversibility*.

**Definition.** The Markov chain  $X_t$  is *reversible* if the probability for  $X_t$  being in state  $x_i$  at time  $t = 0$  and reaching state  $x_j$  at time  $t = s$  is the same as the probability for being in state  $x_j$  at time  $t = 0$  and reaching state  $x_i$  at time  $t = s$  for any  $x_i, x_j \in \mathcal{A}$  and any  $s > 0$ . This requirement is equivalent to the *detailed balance equations*:

$$\begin{aligned} \pi_i P_{ij}(t) &= \pi_j P_{ji}(t), & \forall t > 0 \\ \pi_i Q_{ij} &= \pi_j Q_{ji}. \end{aligned}$$

Now we are able to define an EMP.

**Definition.** We call a time continuous Markov chain  $X_t$  on the set of states  $\mathcal{A}$  an *evolutionary Markov process (EMP)* with the stationary distribution  $\pi$  on the states if and only if:

- $X_t$  is time homogeneous.
- $X_t$  is stationary and the initial distribution  $\pi^{(0)}$  is the stationary distribution  $\pi$ . Therefore  $X_t$  is distributed according to  $\pi$  for all  $t \in \mathbb{R}_0^+$ .
- $X_t$  is irreducible:  $P_{ij}(t) > 0$  for all  $t > 0$  and  $x_i, x_j \in \mathcal{A}$ , i.e.  $\pi$  is unique.
- $X_t$  is calibrated to 1 PEM:  $\sum_i \pi_i Q_{ii} = -0.01$ .
- $X_t$  is reversible:  $\pi_i P_{ij}(t) = \pi_j P_{ji}(t)$  for all  $t > 0$  and all  $x_i, x_j \in \mathcal{A}$ .

## 8.3 Nucleotide Substitution Models

In this section we focus on two important models for nucleotide substitutions, the Jukes-Cantor model and the Kimura 2-parameter model. Both models make assumptions on the rate matrix  $Q$ . Given  $Q$ , the EMP is fully described, as the stationary distribution  $\pi$  can be obtained by solving  $\pi Q = 0$  with  $\sum_{x_i \in \mathcal{A}} \pi_i = 1$ .

### 8.3.1 The Jukes-Cantor Model (JC)

The *Jukes-Cantor model* assumes that each substitution occurs at equal rate  $\alpha$ . Thus the rate matrix is

$$Q = \begin{pmatrix} -3\alpha & \alpha & \alpha & \alpha \\ \alpha & -3\alpha & \alpha & \alpha \\ \alpha & \alpha & -3\alpha & \alpha \\ \alpha & \alpha & \alpha & -3\alpha \end{pmatrix}.$$

The Jukes-Cantor model is reversible and the stationary distribution is the uniform distribution  $\pi = (\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$ . Calibration to 1 PEM yields

$$E = - \sum_{i=1}^4 \pi_i Q_{ii} = 3\alpha \stackrel{!}{=} 0.01,$$

$$\alpha = 1/300.$$

Due to the simple structure of  $Q$ ,  $P(t) = \exp(tQ)$  can be calculated explicitly. The transition probability matrix is

$$P(t) = \begin{pmatrix} 1 - 3a_t & a_t & a_t & a_t \\ a_t & 1 - 3a_t & a_t & a_t \\ a_t & a_t & 1 - 3a_t & a_t \\ a_t & a_t & a_t & 1 - 3a_t \end{pmatrix},$$

where

$$a_t = \frac{1 - \exp(-4\alpha t)}{4} = \frac{1 - \exp(-4t/300)}{4}.$$

**Jukes-Cantor correction.** A simple application of the Jukes-Cantor model is the following:

Given an alignment of two DNA sequences we want to estimate the evolutionary distance between the sequences. Let  $n$  denote the length of the alignment and  $u$  the number of mismatches. A naive distance estimator may take the relative amount of observed substitutions into account only, e.g.  $D = u/n$ . But  $D$  underestimates the amount of substitutions, since multiple and back mutations may have been occurred. Think of a nucleotide, e.g. ‘A’, being substituted by a ‘G’ and the ‘G’ being substituted by an ‘A’ again in the course of evolution. We cannot observe such an event. But we can calculate the probability that any nucleotide remains the same at any site in time  $t$ :

$$\begin{aligned} Pr(X_t = i | X_0 = x_i) &= \sum_{x_i \in \mathcal{A}} \pi_i P_{ii}(t) \\ &= 4 \cdot \frac{1}{4} (1 - 3a_t) \\ &= \frac{1 + 3 \exp(-4\alpha t)}{4} \end{aligned}$$

We also know that  $3\alpha$  substitutions are expected per site and per time unit:  $E = -\sum_i \pi_i Q_{ii} = 4 \cdot (3/4)\alpha = 3\alpha$ . We denote the number of expected substitutions per 100 sites in time  $t$  by  $d$ :  $d = 300\alpha t$  PEM and therefore  $\alpha t = d/300$  PEM. We can replace  $\alpha t$  in the above equation and get

$$Pr(X_t = x_i | X_0 = x_i) = \frac{1 + 3 \exp(-4d/300 \text{ PEM})}{4}.$$

We observe  $u$  mismatches at  $n$  sites and establish

$$Pr(X_t = x_i | X_0 = x_i) = \frac{n - u}{n} = 1 - D = \frac{1 + 3 \exp(-4d/300 \text{ PEM})}{4}.$$

Solving for  $d$  yields

$$d = -\frac{300}{4} \ln \left( 1 - \frac{4}{3} D \right) \text{ PEM}.$$

This is known as *Jukes-Cantor correction* (of the linear distance estimator  $D$ ). If  $d \ll n$ , then  $\ln(1 - (4/3)D) \approx -(4/3)D$  and  $d \approx 100 \cdot D$  PEM. If there are more than five substitutions per 100 sites, that is  $D > 0.05$ , then  $d > 100 \cdot D$  PEM. Note that if  $D > 3/4$ ,  $d$  becomes undefined because the argument of the logarithm becomes negative. Often the score is multiplied by some constant (e.g. 10) and then rounded to the nearest integer in order to facilitate calculations.

### 8.3.2 Kimura 2-Parameter Model (K2P)

Transitions ( $A \leftrightarrow G$  and  $C \leftrightarrow T$ ) are more frequently observed than transversions as  $A$  and  $G$  are purines and  $C$  and  $T$  are pyrimidines. The *Kimura 2-parameter model* takes that into account by introducing different rates for transitions ( $\alpha$ ) and transversions ( $\beta < \alpha$ ). With the order AGCT of nucleotides, the rate matrix  $Q$  is

$$Q = \begin{pmatrix} -\alpha - 2\beta & \alpha & \beta & \beta \\ \alpha & -\alpha - 2\beta & \beta & \beta \\ \beta & \beta & -\alpha - 2\beta & \alpha \\ \beta & \beta & \alpha & -\alpha - 2\beta \end{pmatrix}.$$

For this model, which comprises two parameters, the calculation and calibration of the stationary distribution is analytically feasible. For more parameters, these computations would become very hard. The stationary distribution is also uniform, and calibration to 1 PEM yields  $\alpha + 2\beta = 1/100$ . The probability transition matrix is

$$P(t) = \begin{pmatrix} 1 - (a_t + 2b_t) & a_t & b_t & b_t \\ a_t & 1 - (a_t + 2b_t) & b_t & b_t \\ b_t & b_t & 1 - (a_t + 2b_t) & a_t \\ b_t & b_t & a_t & 1 - (a_t + 2b_t) \end{pmatrix},$$

where

$$\begin{aligned} a_t &:= (2E_t - e_t)/4, & E_t &:= 1 - \exp(-2t(\alpha + \beta)), \\ b_t &:= e_t/4, & e_t &:= 1 - \exp(-4t\beta). \end{aligned}$$

## 8.4 Modeling Amino Acid Replacements

There are only four states in Markov chains describing the evolution of DNA sequences, the four nucleotides. The models presented in the previous sections make simple assumptions on the rate matrices  $Q$ . This kind of procedure is not suited when modeling protein evolution. There are 20 amino acid residues with a variety of replacement frequencies being distributed irregularly. In general one wants to estimate the rate matrix  $Q$  or the probability transition matrix  $P$  by means of amino acid sequence alignments.

### 8.4.1 Parameter Estimation

The strategy of Dayhoff *et al.* [36] is to estimate the 1-step probability transition matrix  $P^{(1)}$  of a time discrete Markov chain and to extrapolate to higher PAM distances. For this purpose pairwise alignments of closely related sequences are collected which are assumed to be correct and which contain approximately 1% mismatches. Entry  $m_{ij}$  of a matrix  $m$  holds the frequency a pair of residues  $(x_i, x_j)$  occurs with, where each pair of residues is counted twice ( $(x_i, x_j)$  and  $(x_j, x_i)$ ) such that the Markov chain is reversible. Variable  $f_i$  denotes the frequency the residue  $x_i$  occurs, and  $N$  is the number of residue pairs. The estimator is based on the equation

$$\frac{m_{ij}}{N} = f_i \cdot P_{ij}^{(1)}.$$

$P^{(1)}$  is calibrated to 1 PAM and the  $k$ -step transition matrices are obtained by  $P^{(k)} = P^k = PP^{k-1}$ , for  $k \geq 2$ .

The main disadvantage of this method is that closely related sequences are taken into account only. It is desirable to exploit sequence alignments of different evolutionary distances for parameter estimation. Say we want to estimate the rate matrix  $Q$  of an EMP. The problem is that the estimator for  $Q$  clearly should account for the evolutionary divergence of each alignment in the dataset. But evolutionary divergence of an alignment has to be estimated by means of the model parameters  $Q$ . Müller and Vingron [35] present an iterative approach cycling between estimating evolutionary distances of sequences in an alignment and updating the current rate matrix  $Q$ .

### 8.4.2 Score Matrices

Searching a protein database with a query protein requires a similarity measure on amino acid sequences. Usually similarity measures on amino acid sequences are based on a similarity measure on the residues which is stored in a score matrix  $(S_{ij})$ .  $S$  assigns a score to each pair of residues  $(x_i, x_j)$ . The PAM matrices are based on a Markov model. The rationale is the following: Similar residues are replaced by each other more frequently than less similar residues. Vice versa one defines similarity by means of replacement frequencies and takes the distribution of residues into account. Consider the following ratio:

$$\frac{\pi_i P_{ij}(t)}{\pi_i \pi_j} = \frac{\pi_j P_{ji}(t)}{\pi_i \pi_j}.$$

The numerator is the probability to observe the pair  $(x_i, x_j)$  in a pair of sequences which have evolved according to the model with the transition matrix  $P(t)$ . The denominator is the

probability to observe the pair  $(x_i, x_j)$  in independent sequences. The *score* is defined as the logarithm of this ratio,

$$S_{ij}(t) := \ln \frac{\pi_i P_{ij}(t)}{\pi_i \pi_j}.$$

The score is positive if the pair  $(x_i, x_j)$  frequently occurs in evolutionarily related sequences, otherwise negative.

### 8.4.3 Maximum Likelihood Estimation of Evolutionary Distances

Given two sequences  $A = (A_1, \dots, A_n)$  and  $B = (B_1, \dots, B_n)$  of the same length which have evolved according to an EMP with rate matrix  $Q$ . We want to estimate the evolutionary distance between  $A$  and  $B$  in PEM units. The probability that  $A$  and  $B$  have evolved from each other in  $t$  time units is

$$Pr(A, B; t) := \prod_{k=1}^n \pi_{A_k} P_{A_k, B_k}(t).$$

We consider  $Pr(A, B; t)$  as likelihood function of time  $t$  to be estimated and try to find a  $t$  such that the probability for the observed data  $A$  and  $B$  is maximal (Maximum Likelihood Estimator). Maximizing the logarithm yields the same  $t$ . Note that the terms  $\pi_{A_k}$  do not depend on  $t$ . Therefore we try to find the maximum of the function

$$\mathcal{L}(t) := \sum_{k=1}^n \log(P_{A_k, B_k}(t)).$$

In general this is done numerically.

# 9 Maximum Likelihood Trees

## 9.1 Computing the Likelihood of a Given Tree

In the context of reconstructing phylogenetic trees from molecular sequence data, the likelihood function which under a fixed model returns for given data the likelihood that these data were produced under the model,  $L(\text{data}) := Pr(\text{data} \mid \text{model})$ , can be written as

$$L(\text{alignment}) = Pr(\text{alignment} \mid \text{tree and evolutionary model}).$$

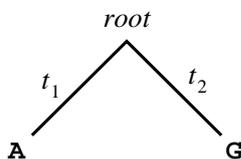
In the following we study a method to reconstruct the most likely tree for a given multiple sequence alignment, under a given evolutionary model.

### Observations:

1. The tree with the highest likelihood can be found by computing the likelihood for each possible tree topology and then choosing the one with the highest likelihood value.
2. The likelihood can be computed for each site of the alignment independently. The total likelihood is then the product over all characters (alignment positions).

Hence, for the moment we will assume that we have given a fixed tree topology with a single character state at each leaf and an evolutionary model that provides us for two character states  $x_i$  and  $x_j$  and an evolutionary distance  $t$  with a probability  $P_{ij}(t)$  that  $x_i$  has changed to  $x_j$  within time  $t$ .

### A simple example:



Assuming the reduced alphabet  $\mathcal{A} = \{A, G\}$ , the total likelihood of this tree is the sum of the likelihoods for the two possible assignments of character states at the root:



$$\begin{aligned}
 L_{\text{total}} &= L_{A \text{ at root}} + L_{G \text{ at root}} \\
 &= P(A \text{ at root})P_{AA}(t_1)P_{AG}(t_2) + P(G \text{ at root})P_{GA}(t_1)P_{GG}(t_2).
 \end{aligned}$$

$P(i \text{ at root})$  is often chosen as the background frequency  $\pi_i$  of state (DNA base)  $x_i$ . Often log-likelihoods are used for easier computation.

After the relation between likelihood and branch lengths is established, the branch lengths  $t_i$  can be adjusted. In this simple example, a short calculation allows to maximize the likelihood using  $\frac{\partial}{\partial t_k} \ln L_{\text{total}} = 0$ . In larger examples, the maximum is usually computed numerically, e.g. by Newton's method.

**The general method:** To compute the likelihood for a larger tree, we use a dynamic programming algorithm. For each vertex  $v$ , we define the *conditional likelihood*  $L_i(v)$  as the likelihood of the subtree below  $v$  given that the character state in  $v$  is  $x_i$ . Then, as above, we have given a tree labelled with one character state at each leaf and an evolutionary model. The algorithm goes as follows.

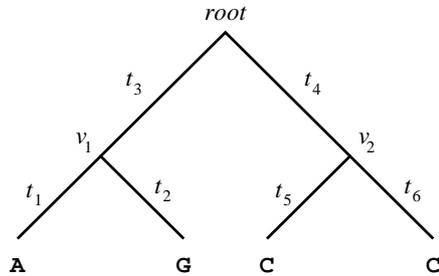
1. Choose an arbitrary root.
2. Traverse the tree bottom-up from the leaves to the root and do the following:
  - a) Assign the likelihood at the leaves: For each leaf  $v$ ,  $L_i(v) = \delta_{ij}$  where  $x_j$  is the character at leaf  $v$ .
  - b) At an internal node  $v$  compute all conditional likelihoods  $L_i(v)$ . Therefore the likelihoods of the different branches are multiplied:

$$L_i(v) = \prod_{v' \text{ child of } v} \sum_{x_j \in \mathcal{A}} P_{ij}(t_{v \rightarrow v'}) L_j(v').$$

3. The total likelihood of the tree is the sum of all conditional likelihoods at the root, weighted by the background probability  $\pi_i$  of the respective character state (DNA base)  $x_i$ :

$$L = \sum_{x_i \in \mathcal{A}} \pi_i L_i(\text{root}).$$

**A larger example:** The following figure shows a rooted tree with four leaves.



Then, for example,

$$L_A(v_1) = P_{AA}(t_1)P_{AG}(t_2)$$

or, more general, for the root

$$L_i(\text{root}) = \left( \sum_j P_{ij}(t_3)L_j(v_1) \right) \left( \sum_j P_{ij}(t_4)L_j(v_2) \right).$$

The total likelihood is:

$$L = \pi_A L_A(\text{root}) + \pi_G L_G(\text{root}) + \pi_C L_C(\text{root}) + \pi_T L_T(\text{root})$$

**Adjusting branch lengths:** An algorithm to adjust branch lengths in the general case is the following.

Repeat for all branches several times:

1. Choose an edge.
2. Choose a node incident to this edge as the root.
3. Compute the maximum likelihood for the rest of the tree as described above.
4. Choose the branch length such that the likelihood of the whole tree is maximized.

**Notes** on the Maximum Likelihood method:

- Changing the model may change the maximum likelihood tree.
- Maximum likelihood is in spirit similar to maximum parsimony, but (1) the cost of a change in parsimony is not a function of the branch length; (2) maximum parsimony only looks at the single, lowest cost solution, whereas maximum likelihood looks at the combined likelihood for all solutions (ancestral states) consistent with the tree and branch lengths.
- The method is extremely time consuming. Heuristic methods exist, e.g. Quartet Puzzling: For each quartet (group of 4 sequences), (1) consider all of the three possible tree topologies and compute their likelihoods, (2) compose intermediate trees from the quartet trees (repeat multiple times), and (3) construct a majority rule consensus tree from all the intermediate trees; then optimize the branch lengths. For details on the Quartet Puzzling method, see Chapter 12.

## 9.2 Consistency

There is an old discussion about which method for reconstructing phylogenetic trees is the best one. Without going too much into the details, here is a list of possible criteria to compare the different methods:

- *Consistency*: The tendency of the method to converge on the correct answer as more data (characters) are added.
- *Efficiency* (or *power*): the ability of the method to recover the correct answer with limited amounts of data.
- *Robustness*: A method is robust if it is relatively insensitive to violations of its assumptions.
- *Computational speed*: The length of time a method takes to solve a problem.
- *Versatility*: What kind of information can be incorporated into the analysis?

In the following we want to focus on consistency, which has received much attention, although it also has its weaknesses. Our discussion of consistency is based on [37] and [1, Chapter 12].

**Definition:** A method is *consistent* if it converges on the correct tree when given infinite data.

All methods are consistent when their assumptions are met, and all methods are inconsistent if their assumptions are sufficiently violated. Therefore, one has to specify the conditions under which a method is consistent.

For example, most distance-based methods (except UPGMA) are consistent under the Jukes-Cantor model. Maximum parsimony can be made consistent by using a Hadamard transformation [38] to correct the data. Maximum likelihood is consistent by definition.

The notion of consistency is not necessarily useful to assess the practical value of a method: A method can be consistent, but very inefficient (like Lake's method of invariants [39, 40]) which means that in practice one will never be able to collect enough data to get a good result with high likelihood.

To test the accuracy of a method in practice, one can apply the method to real (experimentally verified) or numerical (simulated) data. Experimental studies are expensive and time consuming. Simulations are cheap and fast, but might be problematic (model realism, simulation bias).

**Part V**

**Advanced Topics**

# 10 Assessing the Quality of Reconstructed Trees

Once a phylogenetic tree is constructed, it is important to ensure that the result is not the result of some algorithmic artifacts. For example, remember the discussion in Section 7.1.1 about methods that always reconstruct a tree, even if the underlying data is not at all tree-like. In such a case, the result may be some tree, but if the data were looking only slightly different, the tree could be rather different.

One way towards solving this problem is the Splits decomposition method presented in Chapter 7 that produces a tree only for tree-like data, and otherwise a network that shows how strong and in which region the data deviate from a tree.

A different way to assess the reliability of a reconstructed phylogenetic tree is to apply the bootstrap, which is explained in the first part of this chapter. The second part a way to judge about the tree-likeness of data using quartets, i.e. subsets of the data of size four: *likelihood mapping*.

## 10.1 Bootstrapping

The standard reference for the Bootstrap method is the book by Efron and Tibshirani [41]. We first explain the general method, and then how it can be applied to phylogenetic tree methods.

### 10.1.1 The General Idea of Bootstrapping

The general idea of bootstrapping is the following: Let

$$\mathcal{X}_n = (X_1, X_2, \dots, X_n)$$

be a vector of  $n$  samples drawn from an unknown distribution. Let  $T = T(\mathcal{X})$  be some statistic on the sample, e.g. the median. Since the statistic depends on the actual sample (it would be different with a new sample), it can be seen as a random variable, following

the sample statistic. In order to estimate the error of the statistic, we would like to know its standard error. In principle, if the underlying distribution of  $\mathcal{X}$  was known, it would be possible to (a) increase  $n$  or (b) re-draw  $\mathcal{X}$ .

Usually the true distribution is unknown, though, and hence the following trick is used, called *bootstrapping*: Use the sample data  $\mathcal{X}$  and draw from them  $n$  samples with repetition:

$$\mathcal{X}_n^* = (X_1^*, X_2^*, \dots, X_n^*).$$

This can be repeated several times, giving  $B$  *bootstrap replicates*  $\mathcal{X}_n^{*b}$  for  $b = 1, \dots, B$ . For each replicate, the statistic is computed. This allows, e.g., to compute a standard error.

A typical example is a (usually small) set of treatment and control data from a medical experiment (where true replicates are usually very expensive).

### 10.1.2 Bootstrapping in the Phylogenetic Context

In the context of phylogenetic tree reconstruction, bootstrapping was introduced by Felsenstein [42].

Here the sample data are the columns of a multiple alignment. Resampling is done on columns of the alignment. Each replicate is again a multiple alignment of  $n$  columns, but some columns may be duplicated, others may be missing. Each of these alignments gives rise to a tree. The various trees may be different.

Note that bootstrapping is possible for any tree reconstruction method.

Often the bootstrap replicates are used to measure the support for the different parts of the original tree, i.e. the tree  $T$  that was created from the original alignment. A simple way to do this is the following. Each edge in  $T$  defines a split of the taxa. Then one counts for each split (edge) of  $T$  in how many of the bootstrap replicates this split is also realized. This value, often written as a percentage value, is called the *bootstrap support* of the edge.

## 10.2 Likelihood Mapping

Likelihood mapping is a pictorial way to estimate the tree-likeness of aligned molecular sequence data based on likelihood estimation.

**1. Four sequences.** Given four sequences, let  $L_i$  be the maximum likelihood for topology  $i$ . Define  $p_i = L_i / (L_1 + L_2 + L_3)$ . The  $p_i$  define a point in the two-dimensional simplex (for example, illustrated by an equilateral triangle). Each corner corresponds to a topology. See Figures 2 and 3 in [43].

**2. The general case.** There are  $\binom{n}{4}$  quartets. Draw all or a large number of the points. The distribution of points in the resulting picture shows the tree-likeness of the data set. If many points are not in the corners, the data is star-like or not tree-like at all. (Note: The opposite is not generally true.) Figure 4 in [43] shows interesting simulation results for star-like and tree-like data.

# 11 Consensus trees and consensus methods

Given several trees with the same set of elements at their leaves, it is often desirable to compute an average or *consensus* tree of these trees.

## 11.1 The Strict Consensus Tree

The method described here is called the *strict consensus tree*. To this end, one collects all splits that have more than 50% support from the given trees.

**Theorem.** These splits form a tree.

*Proof:*

1. The splits chosen in this way are all pairwise compatible. To see this, assume that there is a pair of splits that is not compatible. Then one can easily derive a contradiction.
2. Remember that for binary characters, pairwise compatibility implies that there exists a perfect phylogeny (2nd version of Gusfield's theorem in Section 3.3), hence they form a (possibly multifurcating) tree.

□

**Definition.** The tree that realizes all these splits with more than 50% support is called the *strict consensus tree*.

Sometimes the strict consensus tree contains highly multifurcating vertices. Therefore it may be further refined by (greedily) adding all splits that fit into the tree and annotating the edges with support values.

## 11.2 Bayesian Methods

A recent and very successful development are Bayesian methods for reconstructing phylogenetic trees. The most popular one is “MrBayes” (<http://mrbayes.csit.fsu.edu/>).

## 12 Parent trees and supertree methods

Given several trees with different (but overlapping) sets of elements at their leaves, it is often desirable to combine them into one larger tree.

### 12.1 Agreement Methods

These are methods that produce a tree that represents an agreement of the structures from the input trees. Examples are: Gordon's strict method, MinCut supertree (including BUILD), strict consensus merger.

### 12.2 Optimization Methods

These are methods that maximize the fit of the constructed tree to the set of source trees according to some objective function.

Examples are: average consensus, Bayesian supertrees, gene tree parsimony, matrix representation methods (MRC, MRF, MRP), quartet supertrees.

### 12.3 One example: Quartet Puzzling

A fast method to compute trees that approximate the maximum likelihood tree is *Quartet Puzzling* [44]. It contains two interesting steps where different trees are compared and joined together.

The procedure is the following:

1. Compute the maximum likelihood solution for all quartets of the data set. (Any other tree reconstruction method for four taxa would work as well.)

2. Piece together the optimal quartets, given an order: one taxon after the other is inserted. Let  $E$  be the next taxon to be inserted. For each edge count the number of quartets containing  $E$  that imply  $E$  to be inserted *not* in this edge. The edge with the lowest overall count wins. (See Figure 2 in [44].)
3. Since the method of Step 2 is order-dependent, it is repeated several times for different input orders. The resulting trees are combined into a so-called strict consensus tree by the method described in Section 11.1.

Quartet puzzling often yields trees rather similar to the maximum likelihood tree. The method, while still slow for many taxa, is much faster than computing the exact maximum likelihood tree. For faster computations, there exists a parallelized version of quartet puzzling [45].

## 13 Fast Convergent Methods

The following is from [46] and [47].

**Definition.** A phylogeny reconstruction method  $M$  is *statistically convergent* under a model of evolution if, for every model tree  $(T, \{\lambda(e)\})$  and every  $\varepsilon > 0$ , there is a sequence length  $k$  such that  $M$  recovers the true tree with probability at least  $1 - \varepsilon$ , when the method is given sequences of length at least  $k$ , generated on the tree  $T$ .

**Definition.** The *convergence rate* of a method  $M$  is the rate at which it converges to 100% accuracy as a function of the sequence length.

**Definition.** Let  $f, g \geq 0$ . Define  $\mathcal{M}_{f,g} = \{(T, \{\lambda(e)\}) : \forall e \in E(T), f \leq \lambda(e) \leq g\}$ .

**Definition.** A phylogenetic reconstruction method  $M$  is (*absolute*) *fast-converging* for the model  $\mathcal{M}$  if, for all positive  $f, g, \varepsilon$ , there is a polynomial  $p$  such that, for all  $(T, \{\lambda(e)\}) \in \mathcal{M}_{f,g}$ , on a set  $S$  of  $n$  sequences of length at least  $p(n)$  generated on  $T$ , we have  $\Pr[M(S) = T] > 1 - \varepsilon$ .

Fast converging methods:

- Harmonic Greedy Triplets and Four Point Condition [48].
- Disk Covering Method with Neighbor Joining [49].

DCM-NJ:

Input:  $n \times n$  distance matrix  $d_{ij}$ .

Algorithm:

1. For each choice of threshold  $q \in \{d_{ij}\}$ , apply disc covering method using NJ, yielding several (not necessary binary) trees, which then by some heuristic method are made binary, giving a collection  $\mathcal{T}$  of binary trees.
2. Select the best tree from  $\mathcal{T}$  by some method: SQS (short quartet support), TS (threshold support), MP (maximum parsimony), ML (maximum likelihood).

# 14 The Method of Invariants

## 14.1 Motivation: Development of a Consistent Method

We have seen that Maximum Parsimony is inconsistent, i.e. even with an infinite amount of data (infinitely long sequences), one will not necessarily obtain the correct tree (remember the discussion about long branch attraction). The reason is a high amount of *homoplasies* (parallel substitutions to the same amino acid) along the long branches that overwhelm the informative character changes along the short internal branch.

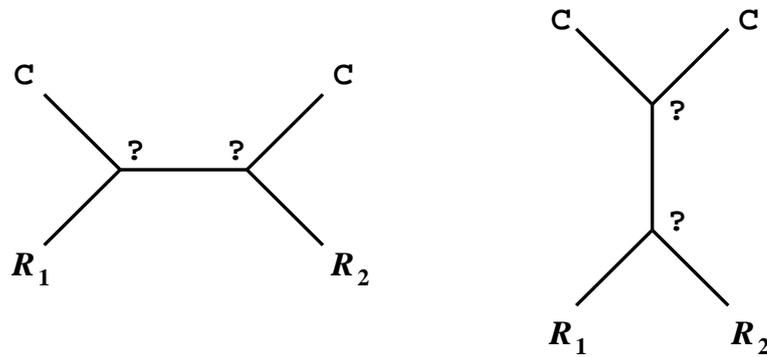
Is it possible to distinguish informative changes from homoplasies? Here are a few possibilities:

- Addition of taxa that lead to branches which split the long edges. (But: taxa must exist, be known, be available; also methods become slower with more taxa.)
- To avoid long branches, use only most conserved sequences (most conserved molecules, and most conserved regions and codons; for example, it is common practice to ignore the third codon position).
- Distance correction in distance based methods.
- Lake's method of invariants, see below.

## 14.2 The Method of Invariants

The *method of invariants* was introduced by Cavender and Felsenstein [50] and by Lake [39] for the case of four taxa. A more general treatment can be found in Evans and Speed [51]. The method is simpler than maximum likelihood in the sense that it is not necessary to estimate parameters describing the evolutionary process. Initially, the method was also called *evolutionary parsimony*.

The method is best described for four-taxon trees (quartets). (For a generalization, see Evans and Speed [51].) We assume that we have given a multiple alignment of four DNA sequences.



**Figure 14.1:** Two trees with different probabilities for  $R_1 = R_2$ .

The method makes use of the fact that the four DNA bases can be divided into two types: A,G are *purines* (R, large) and C,T are *pyrimidines* (Y, small). Substitutions of type  $A \leftrightarrow G$  and  $C \leftrightarrow T$  are *transitions*, all other substitutions are *transversions*. It is generally observed that transitions occur more often than transversions.

The method of invariants is based on three assumptions:

1. Substitutions at a given sequence position are independent.
2. A balance exists among specific classes of transversions (C for example, changes to A or G with equal probabilities).
3. Insertions and deletions can be safely ignored.

**Motivating example.** Compare the two trees in Figure 14.1. In each case, two leaves are labeled with  $C$ . Assume that the other two leaves (labeled  $R_1$  and  $R_2$ ) are (possibly different) purines. Now consider how often one would expect  $R_1$  to be equal to  $R_2$ , and how often the two are expected to be different. It is rather easy to argue intuitively (and a bit more complicated to show formally) that in the right tree the probability that  $R_1 = R_2$  is higher than in the left tree. The method of invariants tries to turn this argument around: By looking at alignment columns with two purines and two pyrimidines, count how often the two purines/pyrimidines are the same and how often they are different; then infer from this the correct tree topology.

**More general example.** Given a multiple alignment of four sequences, assume we extract only those columns where two specific sequences have bases of one type (both purine or both pyrimidine) and the other two sequences have bases of the opposite type. For example,

in the following alignment, in each column, rows 1 and 2 have the same type and rows 3 and 4 have the opposite type.

```

1: A A G A C T T G
2: A G G G C C T G
3: C T C C G A A T
4: C T T T A G A T

```

We use the following notation: Assume, as in the above example, we have an alignment extraction where rows 1 and 2 are of identical type and rows 3 and 4 are of identical type. Then  $f_{12||34}$  is defined as the number of columns where the bases in rows 1 and 2 are the same, and the bases in rows 3 and 4 are the same. Similarly,  $f_{1|2||34}$  is the number of columns where the bases in rows 1 and 2 differ, but the bases in rows 3 and 4 are identical;  $f_{12||3|4}$  is the number of columns where rows 3 and 4 differ and rows 1 and 2 are identical; and finally  $f_{1|2||3|4}$  is the number of columns where all four bases are different.

For example, in the above alignment we have

$$\begin{aligned}
 f_{12||34} &= 3 \quad (\text{columns 1, 7 and 8}), \\
 f_{1|2||34} &= 1 \quad (\text{column 2}), \\
 f_{12||3|4} &= 2 \quad (\text{columns 3 and 5}), \\
 f_{1|2||3|4} &= 2 \quad (\text{columns 4 and 6}).
 \end{aligned}$$

In a similar way one defines  $f_{13||24}$  etc. for alignment extractions where rows 1 and 3 resp. rows 2 and 4 are of the same type, and  $f_{14||23}$  etc. for alignment extractions where rows 1 and 4 resp. rows 2 and 3 are of the same type.

Finally, define the *net support* (the *invariant*) for topology  $\frac{1}{2} \succ \prec \frac{3}{4}$  by

$$F \left( \frac{1}{2} \succ \prec \frac{3}{4} \right) := f_{12||34} - f_{1|2||34} - f_{12||3|4} + f_{1|2||3|4}.$$

Similar, define

$$F \left( \frac{1}{3} \succ \prec \frac{2}{4} \right) := f_{13||24} - f_{1|3||24} - f_{13||2|4} + f_{1|3||2|4},$$

$$F \left( \frac{1}{4} \succ \prec \frac{2}{3} \right) := f_{14||23} - f_{1|4||23} - f_{14||2|3} + f_{1|4||2|3}.$$

One can make the following observation:

- $F \left( \begin{smallmatrix} 1 \\ 2 \end{smallmatrix} \succ \left\langle \begin{smallmatrix} 3 \\ 4 \end{smallmatrix} \right. \right) \neq 0$  if topology  $\begin{smallmatrix} 1 \\ 2 \end{smallmatrix} \succ \left\langle \begin{smallmatrix} 3 \\ 4 \end{smallmatrix} \right.$  is correct.
- $F \left( \begin{smallmatrix} 1 \\ 2 \end{smallmatrix} \succ \left\langle \begin{smallmatrix} 3 \\ 4 \end{smallmatrix} \right. \right) \approx 0$  if topology  $\begin{smallmatrix} 1 \\ 2 \end{smallmatrix} \succ \left\langle \begin{smallmatrix} 3 \\ 4 \end{smallmatrix} \right.$  is wrong.

To see why this is so, we give a motivation for the four terms. Therefore assume in the following that the true topology is  $\begin{smallmatrix} 1 \\ 2 \end{smallmatrix} \succ \left\langle \begin{smallmatrix} 3 \\ 4 \end{smallmatrix} \right.$ .

**Parsimony term  $f_{12||34}$ :** for small leaf edges, the peripheral branches term and the compensatory term will be small  $\Rightarrow$  same as parsimony.

**Two peripheral branches terms  $f_{12||3|4}$  and  $f_{1|2||34}$ .** two long (misleading) leaf edges will increase the parsimony term of another topology, but also one of the peripheral branches terms of that topology will increase:

- 1,3 long – in  $F \left( \begin{smallmatrix} 1 \\ 3 \end{smallmatrix} \succ \left\langle \begin{smallmatrix} 2 \\ 4 \end{smallmatrix} \right. \right)$ ,  $f_{13||24}$  will increase, but so will  $f_{1|3||24}$ .
- 2,4 long – in  $F \left( \begin{smallmatrix} 1 \\ 3 \end{smallmatrix} \succ \left\langle \begin{smallmatrix} 2 \\ 4 \end{smallmatrix} \right. \right)$ ,  $f_{13||24}$  will increase, but so will  $f_{13||2|4}$ .
- 1,4 long – in  $F \left( \begin{smallmatrix} 1 \\ 4 \end{smallmatrix} \succ \left\langle \begin{smallmatrix} 2 \\ 3 \end{smallmatrix} \right. \right)$ ,  $f_{14||23}$  will increase, but so will  $f_{1|4||23}$ .
- 2,3 long – in  $F \left( \begin{smallmatrix} 1 \\ 4 \end{smallmatrix} \succ \left\langle \begin{smallmatrix} 2 \\ 3 \end{smallmatrix} \right. \right)$ ,  $f_{14||23}$  will increase, but so will  $f_{14||2|3}$ .

**Compensatory term  $f_{1|2||3|4}$ .** frequent transitions (noise) decrease  $f_{12||34}$  and create wrong  $12||3|4$  and  $1|2||34$ . Therefore, the author has introduced the compensating term  $1|2||3|4$  that acts like a transition normalizer.

**Algorithm:**

1. Choose a quartet of aligned sequences 1, 2, 3, 4.
2. Find the alignment positions in which two sequences have purines and two have pyrimidines.
3. Compute the net supports (invariants) for the three possible topologies  $F \left( \begin{smallmatrix} 1 \\ 2 \end{smallmatrix} \succ \left\langle \begin{smallmatrix} 3 \\ 4 \end{smallmatrix} \right. \right)$ ,  $F \left( \begin{smallmatrix} 1 \\ 3 \end{smallmatrix} \succ \left\langle \begin{smallmatrix} 2 \\ 4 \end{smallmatrix} \right. \right)$ , and  $F \left( \begin{smallmatrix} 1 \\ 4 \end{smallmatrix} \succ \left\langle \begin{smallmatrix} 2 \\ 3 \end{smallmatrix} \right. \right)$ . The support for two of the three topologies should be near zero, while the third topology may (or may not) be supported by a non-zero score.
4. (This support may be quantified by a statistical analysis, e.g. the  $\chi^2$  test. But be careful: Often the counts are too low such that the requirements for an applicability of the  $\chi^2$ -test are not fulfilled!)

**General scheme for an invariant:** A function  $F$  of the  $f_{..||..}$  is an *invariant* (for some model of evolution) if  $F$  is zero for all possible choices of model parameters under a (wrong) phylogeny, and  $F$  is typically nonzero for the correct phylogeny.

**Remarks:**

- It is possible that  $F$  is negative. It has been suggested in the literature that this could be interpreted as a result of evolutionary pressure. But the arguments are not very strong.
- One hope for the method of invariants: Assumption of identical distribution for different sites can be weakened more easily than for the maximum likelihood method.

## Bibliography

- [1] D. M. Hillis, C. Moritz, and B. K. Mable, editors. *Molecular Systematics*. Sinauer, Sunderland, MA, 2nd edition, 1996.
- [2] D. Graur and W.-H. Li, editors. *Fundamentals of Molecular Evolution*. Sinauer, Sunderland, MA, 2nd edition, 2000.
- [3] D. Gusfield. *Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology*. Cambridge University Press, New York, NY, 1997.
- [4] R. D. M. Page and E. C. Holmes, editors. *Molecular Evolution: A Phylogenetic Approach*. Blackwell Science, Oxford, England, 1998.
- [5] J. Felsenstein, editor. *Inferring Phylogenies*. Sinauer, Sunderland, MA, 2004.
- [6] Alfred Russel Wallace. *Darwinism. An Exposition of the Theory of Natural Selection with Some of its Applications*. John Wilson & Son, 1889.
- [7] A. J. Griffiths, J. H. Miller, D. T. Suzuki, R. C. Lewontin, and W. M. Gelbart. *An Introduction to Genetic Analysis (6th ed.)*. W.H. Freeman and Company, New York, 1996.
- [8] Dan Gusfield. Efficient algorithms for inferring algorithms for evolutionary trees. *Networks*, 21:19–28, 1991.
- [9] Fumei Lam, Daniel Gusfield, and Srinath Sridhar. Generalizing the four gamete condition and splits equivalence theorem: Perfect phylogeny on three state characters. In *Proc. of the 9th International Workshop on Algorithms in Bioinformatics, WABI 2009*, volume 5724 of *LNBI*, pages 206–219, Berlin, 2009. Springer Verlag.
- [10] Tandy Warnow. Tree compatibility and inferring evolutionary history. *Journal of Algorithmics*, 16:388–407, 1994.
- [11] M. S. Waterman. *Introduction to Computational Biology. Maps, Sequences and Genomes*. Chapman & Hall, London, 1995.
- [12] J. C. Setubal and J. M. Meidanis. *Introduction to Computational Molecular Biology*. PWS Publishing, Boston, MA, 1997.

- 
- [13] W. M. Fitch. Toward defining the course of evolution: Minimum change for a specific tree topology. *SystZool*, 20(4):406–416, 1971.
- [14] J. A. Hartigan. Minimum mutation fits to a given tree. *Biometrics*, 29(1):53–65, 1973.
- [15] D. Sankoff and P. Rousseau. Locating the vertices of a steiner tree in an arbitrary metric space. 9:240–246, 1975.
- [16] D. Sankoff and R. J. Cedergren. Simultaneous comparison of three or more sequences related by a tree. In D. Sankoff and J. B. Kruskal, editors, *Time Warps, String Edits, and Macromolecules: The Theory and Practice of Sequence Comparison*, chapter 9, pages 253–263. Addison-Wesley, Reading, MA, 1983.
- [17] D. Sankoff. Minimal mutation trees of sequences. *SIAM J. Appl. Math.*, 28:35–42, 1975.
- [18] M. D Hendy and D. Penny. Branch and bound algorithms to determine minimal evolutionary trees. *Math. Biosci.*, 59:277–290, 1982.
- [19] J. Felsenstein. *Inferring Phylogenies*. Sinauer Associates, Inc., Sunderland, Massachusetts, 2004.
- [20] L. R. Foulds and R. L. Graham. The Steiner problem in phylogeny is NPcomplete. *Adv. Appl. Math.*, 3:43–49, 1982.
- [21] T. H. Cormen, C. E. Leiserson, and R. L. Rivest. *Introduction to Algorithms*. MIT Press, Cambridge, MA, 1989.
- [22] M. Vingron and A. von Haeseler. Towards integration of multiple alignment and phylogenetic tree construction. *J. Comp. Biol.*, 4:23–34, 1997.
- [23] B. Schwikowski and M. Vingron. The deferred path heuristic for the generalized tree alignment problem. *J. Comp. Biol.*, 4:415–431, 1997.
- [24] P. Bunemann. The recovery of trees from measures of dissimilarity. In J. Hodson et al., editor, *Mathematics and the Archeological and Historical Sciences*, pages 387–359. Edinbhurg University Press, 1971.
- [25] C. D. Michener and R. R. Sokal. A quantitative approach to a problem in classification. *Evolution*, 11:130–162, 1957.
- [26] M. S. Waterman, T. F. Smith, M. Singh, and W. A. Beyer. Additive evolutionary trees. *J. Theor. Biol.*, 64:199–213, 1977.
- [27] W.M. Fitch and E. Margoliash. Construction of phylogenetic trees. *Science*, 155:279–284, 1967.
- [28] A. Rzhetsky and M. Nei. A simple method for estimating and testing minimum-evolution trees. *Mol. Biol. Evol.*, 9:945–967, 1992.

- [29] R. Desper and O. Gascuel. Fast and accurate phylogeny reconstruction algorithms based on the minimum-evolution principle. In R. Guigó and D. Gusfield, editors, *Proceedings of the Second International Workshop on Algorithms in Bioinformatics, WABI 02*, volume 2452 of *LNCS*, pages 357–374, Berlin, 2002. Springer Verlag.
- [30] J.A. Studier and K.J. Keppler. A note on the neighbor-joining algorithm of saitou and nei. *Molecular Biology and Evolution*, 5:729–731, 1988.
- [31] K. Atteson. The performance of the neighbor joining method of phylogeny reconstruction. In B. Mirkin et al., editor, *Mathematical Hierarchies and Biology*, pages 133–148. American Mathematical Society, 1997.
- [32] H.-J. Bandelt and A. W. M. Dress. A canonical decomposition theory for metrics on a finite set. *Adv. Math.*, 92:47–105, 1992.
- [33] H.-J. Bandelt and A. W. M. Dress. Split decomposition: A new and useful approach to phylogenetic analysis of distance data. *Mol. Phylogenet. Evol.*, 1:242–252, 1992.
- [34] D. Bryant and V. Moulton. NeighborNet: An agglomerative method for the construction of planar phylogenetic networks. In R. Guigó and D. Gusfield, editors, *Proceedings of the Second International Workshop on Algorithms in Bioinformatics, WABI 02*, volume 2452 of *LNCS*, pages 375–391, Berlin, 2002. Springer Verlag.
- [35] T. Müller and M. Vingron. Modeling amino acid replacement. *J. Comp. Biol*, 6:761–776, 2000.
- [36] M. O. Dayhoff, R. M. Schwartz, and B. C. Orcutt. A model of evolutionary change in proteins. In M. O. Dayhoff, editor, *Atlas of Protein Sequence and Structure*, volume 5, suppl. 3, pages 345–352. National Biomedical Research Foundation, Washington, D.C., 1979.
- [37] J. P. Huelsenbeck. Performance of phylogenetic methods in simulation. *Syst. Biol.*, 44(1):17–48, 1995.
- [38] M. D. Hendy and D. Penny. Spectral analysis of phylogenetic data. *J. Classif.*, 10:5–24, 1993.
- [39] J. A. Lake. A rate-independent technique for analysis of nucleic acid sequences: Evolutionary parsimony. *Mol. Biol. Evol.*, 4(2):167–191, 1987.
- [40] J. A. Lake. Origin of the eukaryotic nucleus determined by rate-invariant analysis of rRNA sequences. *Nature*, 331:184–186, 1988.
- [41] B. Efron and R. J. Tibshirani. *An Introduction to the Bootstrap*. CRC Press, Boca Raton, FL, 1998.

- 
- [42] J. Felsenstein. Confidence limits on phylogenies: An approach using the bootstrap. *Evolution*, 39(4):783–791, 1985.
- [43] K. Strimmer and A. von Haeseler. Likelihood-mapping: A simple method to visualize phylogenetic content of a sequence alignment. *Proc. Natl. Acad. Sci. USA*, 94:6815–6819, 1997.
- [44] K. Strimmer and A. von Haeseler. Quartet puzzling: A quartet maximum-likelihood method for reconstructing tree topologies. *Mol. Biol. Evol.*, 13(7):964–969, 1996.
- [45] H. A. Schmidt, K. Strimmer, M. Vingron, and A. von Haeseler. TREE-PUZZLE: Maximum likelihood phylogenetic analysis using quartets and parallel computing. *Bioinformatics*, 18(3):502–504, 2002.
- [46] L. Nakhleh, U. Roshan, K. St. John, J. Sun, and T. Warnow. Designing fast converging phylogenetic methods. *Bioinformatics*, 17(Suppl. 1):S190–S198, 2001. (Proceedings of ISMB 2001).
- [47] B. M. E. Moret, U. Roshan, and T. Warnow. Sequence-length requirements for phylogenetic methods. In R. Guigó and D. Gusfield, editors, *Proceedings of the Second International Workshop on Algorithms in Bioinformatics, WABI 02*, volume 2452 of *LNCS*, pages 343–356, Berlin, 2002. Springer Verlag.
- [48] M. Csűrös and M.-Y. Kao. Provably fast and accurate recovery of evolutionary trees through harmonic greedy triplets. *SIAM J. Computing*, 31(1):306–322, 2001.
- [49] D. H. Huson, S. Nettles, and T. J. Warnow. Disk-covering, a fast-converging method for phylogenetic tree reconstruction. *J. Comp. Biol.*, 6(3-4):369–386, 1999.
- [50] J. A. Cavender and J. Felsenstein. Invariants of phylogenies: Simple cases with discrete states. *J. Classif.*, 4:57–71, 1987.
- [51] S. N. Evans and T. P. Speed. Invariants of some probability models used in phylogenetic inference. *Ann. Stat.*, 21(1):355–377, 1993.