

# Handout for the Phylogenetics Lecture

Dirk Metzler

Winter Semester 2012/2013, last updated on 4. February 2013

## Contents

<b>1</b>	<b>Intro: Outline and Tree Notation</b>	<b>2</b>
<b>2</b>	<b>Distance-Based Phylogeny Reconstruction</b>	<b>5</b>
2.1	What is a distance? . . . . .	5
2.2	UPGMA . . . . .	5
2.3	Neighbor Joining . . . . .	7
<b>3</b>	<b>Parsimony in phylogeny reconstruction</b>	<b>8</b>
3.1	Parsimony of a tree . . . . .	8
3.2	Finding parsimonious trees for given data . . . . .	10
3.3	Excursus: measuring the difference between two trees . . . . .	13
3.4	Limitations of the parsimony principle . . . . .	13
<b>4</b>	<b>Maximum-Likelihood (ML) in phylogeny estimation</b>	<b>14</b>
4.1	What is a likelihood? . . . . .	14
4.2	How to compute the likelihood of a tree . . . . .	15
4.3	How to search for the ML tree . . . . .	17
4.4	Consistency of the Maximum-Likelihood method . . . . .	19
4.5	Maximum Parsimony from a probabilistic perspective . . . . .	19
4.6	Maximum likelihood for pairwise distances . . . . .	20
<b>5</b>	<b>Modelling the substitution process on sequences</b>	<b>20</b>
5.1	Transition matrix and rate matrix . . . . .	20
5.2	Residence time . . . . .	22
5.3	Computing $S(t)$ from the rate matrix $R$ . . . . .	23
5.4	A model for transition probabilities in closed form . . . . .	25
5.5	Overview of DNA substitution models . . . . .	25
5.6	Position-dependent mutation rates . . . . .	26
5.7	Convergence into equilibrium . . . . .	28
<b>6</b>	<b>Bayesian phylogeny reconstruction and MCMC</b>	<b>29</b>
6.1	Principles of Bayesian statistics . . . . .	29
6.2	MCMC sampling . . . . .	30
6.3	Checking convergence of MCMC . . . . .	32
6.4	Interpretation of posterior probabilities and robustness . . . . .	33
<b>7</b>	<b>Bootstrapping</b>	<b>34</b>
7.1	The concept of bootstrapping . . . . .	34
7.2	Bootstrap for phylogenetic trees . . . . .	35
7.3	How can we interpret the bootstrap values? . . . . .	36

<b>8</b>	<b>Tests for trees and branches</b>	<b>37</b>
8.1	The Kishino-Hasegawa (KH) test . . . . .	37
8.2	The Shimodaira-Hasegawa (SH) test . . . . .	37
8.3	The SOWH test . . . . .	38
8.4	An approximate Likelihood-Ratio Test (aLRT) . . . . .	38
<b>9</b>	<b>Model selection</b>	<b>39</b>
9.1	Concepts: AIC, hLRT, BIC, DT, Model averaging, and bootstrap again . . . . .	39
9.2	Does model selection matter? . . . . .	42
<b>10</b>	<b>Insertion-Deletion Models for Statistical Alignment</b>	<b>43</b>
10.1	Alignment sampling with pairHMMs . . . . .	43
10.2	Insertions and deletions of more than one site . . . . .	50
10.3	Multiple Alignments . . . . .	52
10.4	Software for joint estimation of phylogenies and alignments . . . . .	55
<b>11</b>	<b>Quantitative Characters and Independent Contrasts</b>	<b>56</b>
11.1	Brownian motions along the branches of the tree . . . . .	56
11.2	Excursus: Multidimensional Normal Distribution . . . . .	57
11.3	Why to use REML . . . . .	60
11.4	Computing Independent Contrasts by Pruning the Tree . . . . .	61
11.5	Software . . . . .	62

# 1 Intro: Outline and Tree Notation

## Tentative plan for the rest of the semester

- Maximum Likelihood v. Maximum Parsimony vs. Distance-based Phylogeny Inference
- Sequence Evolution Models: JC, F81, HKY, F84, GTR, PAM and  $\Gamma$ -distributed rates
- Bootstrap
- MCMC and Bayesian Inference
- Relaxed Molecular Clock and Time Calibration
- Independent Contrasts for Quantitative Traits
- Statistical Alignment (TKF91, TKF92, pairHMMs, multiple HMMs)
- Genetree-Speciestree Reconciliation

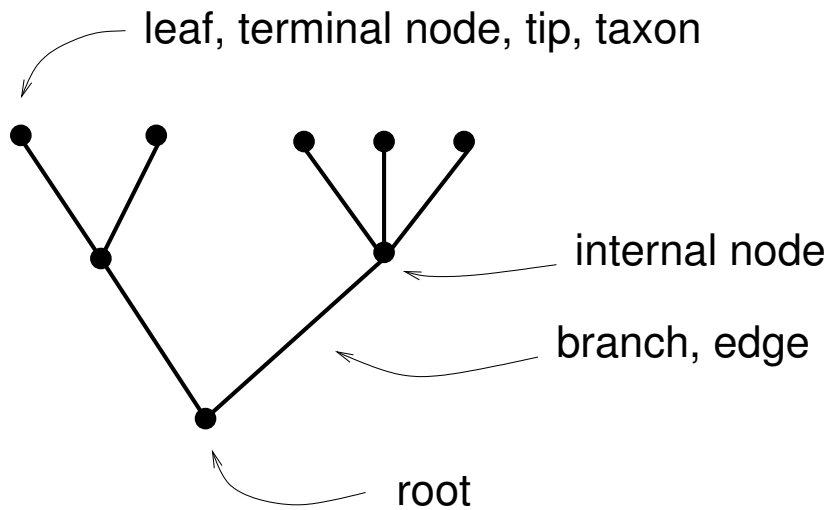
## Aims

- Understand principles and rationales underlying the methods
- Explore available software
- What is efficiently doable, what is difficult?
- What are the strengths and weaknesses of the methods?
- Which method is appropriate for which dataset?
- Learn what is necessary to read papers about new computational methods
- Future directions of phylogenetics

## Recommended Books

## References

- [Fel04] J. Felsenstein (2004) *Inferring Phylogenies*
- [Yang06] Z. Yang (2006) *Computational Molecular Evolution*
- [Niel05] R. Nielsen, [Ed.] (2005) *Statistical Methods in Molecular Evolution*
- [DEKM98] R. Durbin, S. Eddy, A. Krogh, G. Mitchison (1998) *Biological Sequence Analysis*
- [EG05] W. Ewens, G. Grant (2005) *Statistical Methods in Bioinformatics*



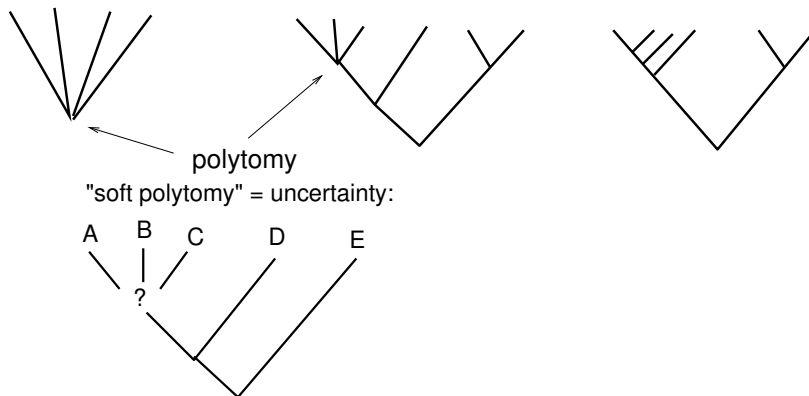
degree of a node = number of edges adjacent to the node

binary tree = fully resolved tree: root has degree two, all other nodes have degree 3

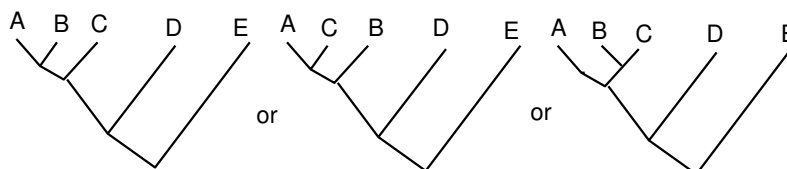
star-shaped tree

partially resolved

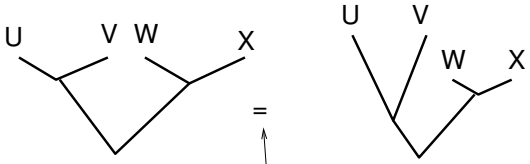
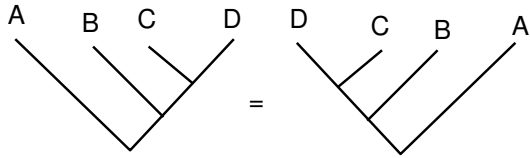
fully resolved



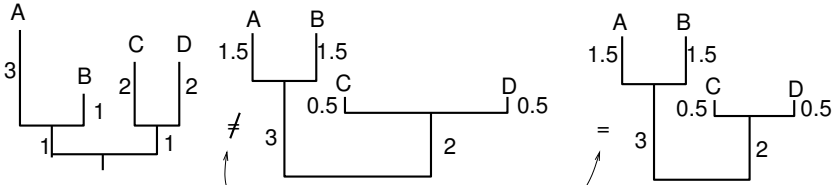
could be



cladogram: branch lengths not meaningful

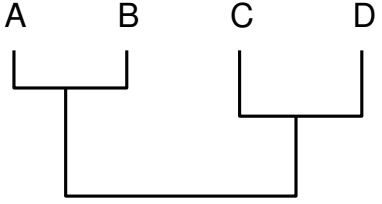


as cladograms

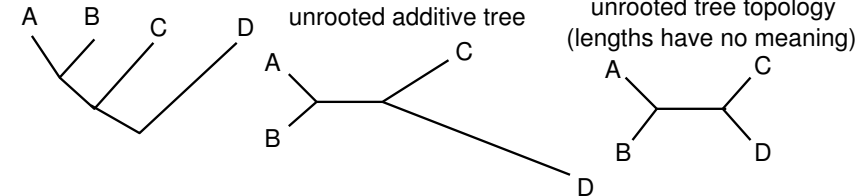


as additive trees

dendrogram = chronogram = ultrametric tree = additive trees that are compatible with molecular-clock assumption, i.e. all tips have the same distance to the root

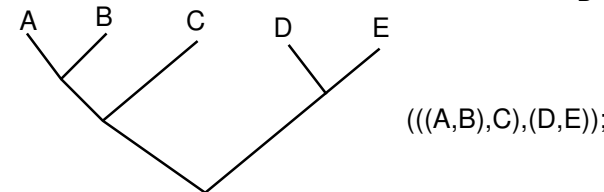


rooted additive tree

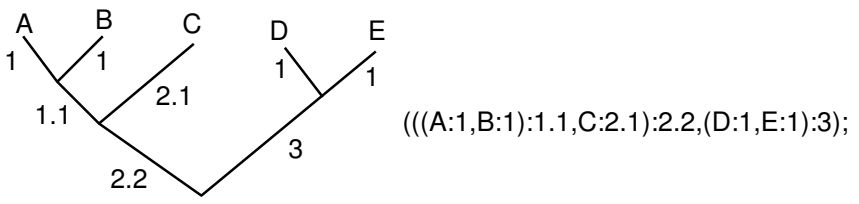


unrooted additive tree

unrooted tree topology (lengths have no meaning)



$((A,B),C),(D,E);$

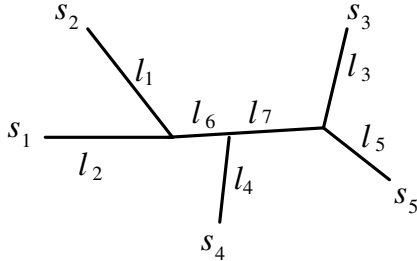


$((A:1,B:1):1.1,C:2.1):2.2,(D:1,E:1):3);$

## 2 Distance-Based Phylogeny Reconstruction

### 2.1 What is a distance?

Given a set of taxa  $S = \{s_1, s_2, \dots, s_n\}$  and matrix of distances  $(d_{ij})_{i,j \leq n}$ , where  $d_{ij}$  is the (estimated) distance between  $s_i$  and  $s_j$  we search for a tree whose tips are labeled with  $S$  and whose edges are labeled with lengths, such that the distances between tips labeled with  $s_i$  and  $s_j$  should be (approximately)  $d_{ij}$  for all  $i, j$ .



For example,  $l_2 + l_6 + l_7 + l_3$  should be (as close as possible to)  $d_{1,3}$

Distances should be *additive*. E.g. if we define the distance to be the number of observed differences between DNA sequences, this distance will not be additive. A more useful distance is the *expected* number of mutations according to a sequence evolution model; more about this later.

To be a proper *distance matrix*,  $(d_{ij})_{i,j \leq n}$  must fulfill the following requirements for all  $i, j, k$ :

- $d_{ij} = d_{ji}$
- $d_{ij} = 0 \Leftrightarrow i = j$
- $d_{ij} + d_{jk} \geq d_{ik}$  (triangle inequality)

### 2.2 UPGMA

UPGMA (**U**nweighted **P**airwise **G**rouping **M**ethod with **A**rithmetic mean, Sokal & Michener, 1985) is hierarchical cluster method using means:

for  $i \leq n$  set  $C_i := \{s_i\}$

$\mathcal{C} := \{C_1, \dots, C_n\}$  is the current set of clusters

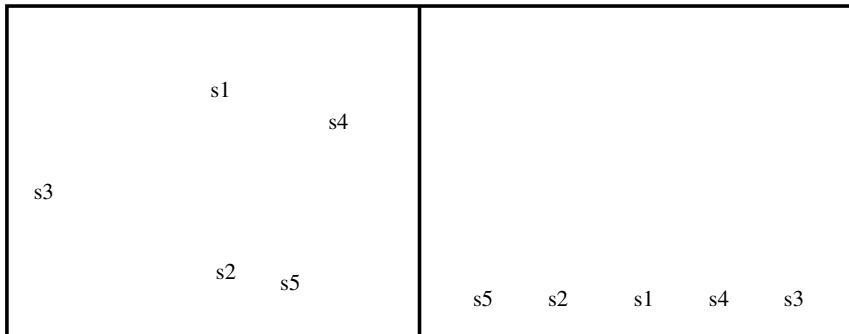
$m := n$

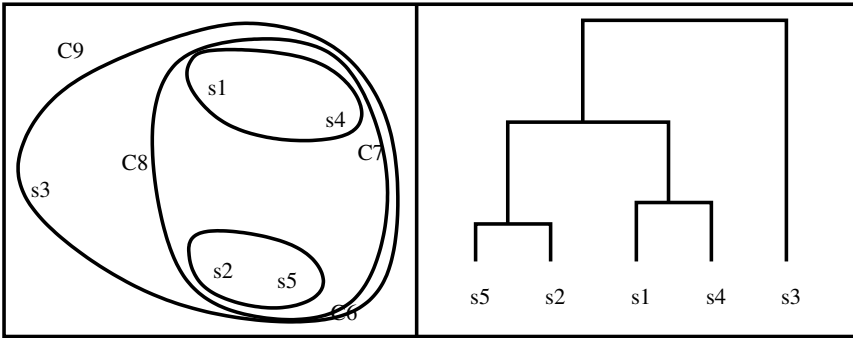
repeat ...

- $m := m - 1$
- find  $C_i, C_j \in \mathcal{C}$  with minimum  $d_{ij} > 0$
- $C_m := C_i \cup C_j$
- $\mathcal{C} := \mathcal{C} \cup \{C_m\} \setminus \{C_i, C_j\}$
- For all  $C_k \in \mathcal{C}$  set

$$d_{km} := d_{mk} := \frac{1}{|C_k| \cdot |C_m|} \sum_{s_x \in C_k, s_y \in C_m} d_{xy}$$

... until  $C_m = \{s_1, \dots, s_n\}$  and  $\mathcal{C} = C_m$ .





common ways to define the distance between clusters  $C$  and  $C'$  cluster algorithms:

**single linkage:**  $d(C, C') = \min_{s_i \in C, s_j \in C'} d_{ij}$

**complete linkage:**  $d(C, C') = \max_{s_i \in C, s_j \in C'} d_{ij}$

**means (like in UPGMA):**  $d(C, C') = \frac{1}{|C| \cdot |C'|} \sum_{i \in C, j \in C'} d_{ij}$

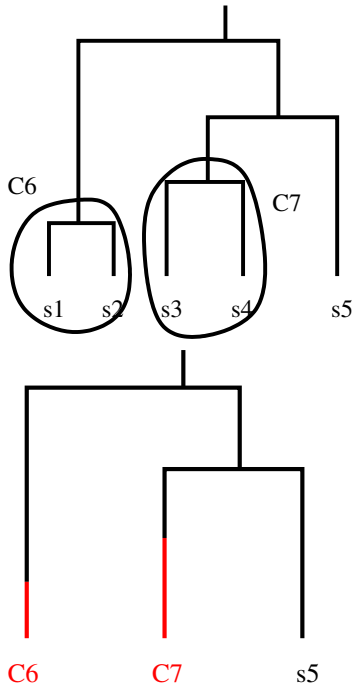
**UPGMA works under ideal conditions**

Assume there is an **ultrametric** tree (i.e. molecular-clock) in which the tips have **exactly** the given distances  $d_{ij}$ . Then, UPGMA will find this tree.

Reason: in the first step UPGMA will correctly join the closest relatives.

As a consequence of the molecular clock assumption, UPGMA will define reasonable distances between the clusters.

Example:



From

$$d_{13} = d_{14} = d_{23} = d_{24}$$

follows

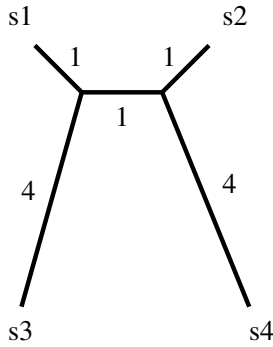
$$d_{67} = \frac{1}{2 \cdot 2} \cdot (d_{13} + d_{14} + d_{23} + d_{24}) = d_{13}$$

This means that we are in the same situation as in the first step: The clusters are tips of an ultrametric tree, and the distances for the clusters are just like the distances of any taxa in the clusters.

Thus, UPGMA will not only get the first step right but also any other step.

**When UPGMA fails**

If the tree is not compatible with molecular-clock assumptions, UPGMA may fail even if the precise distances are known.



In this example, UPGMA will first join  $s_1$  and  $s_2$  and will not have a chance to correct this in any later step.

### Ultrametric distances

**Theorem 1** Let  $D = (d_{ij})_{ij}$  be a distance matrix for  $(s_1, \dots, s_n)$ . The following two properties are equivalent:

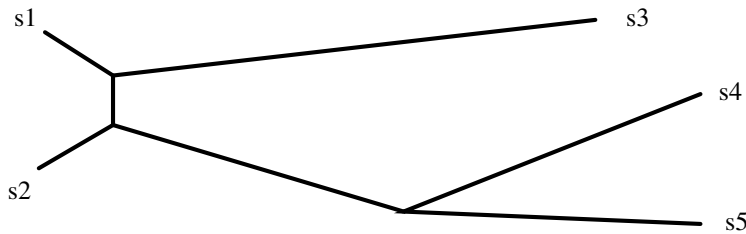
- (a) A binary tree exists that fulfills the molecular-clock assumption and the tips of this tree have the distances given in  $D$ . (The distance between two tips is the sum of the lengths of the edges between them.)
- (b)  $D$  is ultrametric, i.e.

$$\forall \text{sets of three indices } \{i, j, k\} \exists i \in \{i, j, k\} : d_{jk} < d_{ij} = d_{ik}$$

### 2.3 Neighbor Joining

Idea: use modified distances that take into account how far a taxon is to all other taxa

$$D_{ij} := d_{ij} - (r_i + r_j), \quad \text{where } r_i = \frac{1}{n-2} \sum_k d_{ik} = \frac{n-1}{n-2} \cdot \bar{d}_i.$$

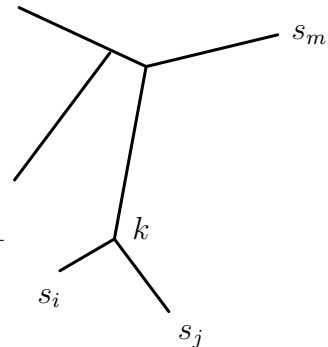


#### Neighbor Joining algorithm (Saitou, Nei, 1987)

Input  $T = \{s_1, \dots, s_n\}$  with distance matrix  $(d_{ij})_{i,j \leq n}$

NeighborJoining( $T$ ):

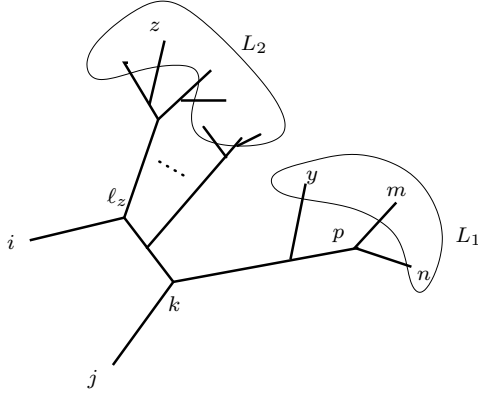
- done if  $n = 1$
- compute all  $D_{ij}$
- find taxa  $s_i$  and  $s_j$  in  $T$  with minimum  $D_{ij}$
- define internal node  $k$  with distances  $\forall_m : d_{km} := \frac{1}{2}(d_{im} + d_{jm} - d_{ij})$
- NeighborJoining( $\{k\} \cup T \setminus \{s_i, s_j\}$ )



**Theorem 2 (Neighbor-Joining theorem, Studier & Keppler, 1988)** If a tree exists whose tips have precisely the distances given by  $(d_{ij})_{ij}$ , then Neighbor-Joining will find this tree.

Sketch of proof: assume that  $i$  and  $j$  are not neighbors and show that  $D_{ij}$  can then not be minimal. Let set of tips  $L_1$  and  $L_2$  be defined as below and assume w.l.o.g.  $|L_1| \leq |L_2|$ . By definition,

$$D_{ij} - D_{mn} = d_{ij} - d_{mn} - \frac{1}{n-2} \left( \sum_u d_{iu} + d_{ju} - d_{mu} - d_{nu} \right).$$



Using additivity, we can show that

$$\begin{aligned} d_{iy} + d_{jy} - d_{my} - d_{ny} \\ = d_{ij} + 2d_{ky} - 2d_{py} - d_{mn} \end{aligned}$$

and

$$\begin{aligned} d_{iz} + d_{jz} - d_{mz} - d_{nz} \\ = d_{ij} - d_{mn} - 2d_{pk} - 2d_{\ell_z k}. \end{aligned}$$

hold for all tips  $y \in L_1$  and  $z \in L_2$ .

We assume that  $|L_1| \geq 2$  and leave the case  $|L_1| = 1$  as an exercise.

It follows from the equations above that

$$\begin{aligned} D_{ij} - D_{mn} &= \frac{\left( \sum_{y \in L_1} 2d_{py} - 2d_{ky} \right) + \left( \sum_{z \in L_2} 2d_{pk} + 2d_{\ell_z k} \right) + 4d_{kp} + 2d_{mn}}{n-2} \\ &> 2d_{pk}(|L_2| - |L_1|)/(n-2) \quad (\text{because of } d_{py} - d_{ky} > -d_{pk}) \\ &\geq 0 \end{aligned}$$

and thus  $D_{ij} > D_{mn}$ , q.e.d.

### 3 Parsimony in phylogeny reconstruction

#### 3.1 Parsimony of a tree

Given  $n$  homologous DNA or protein sequences

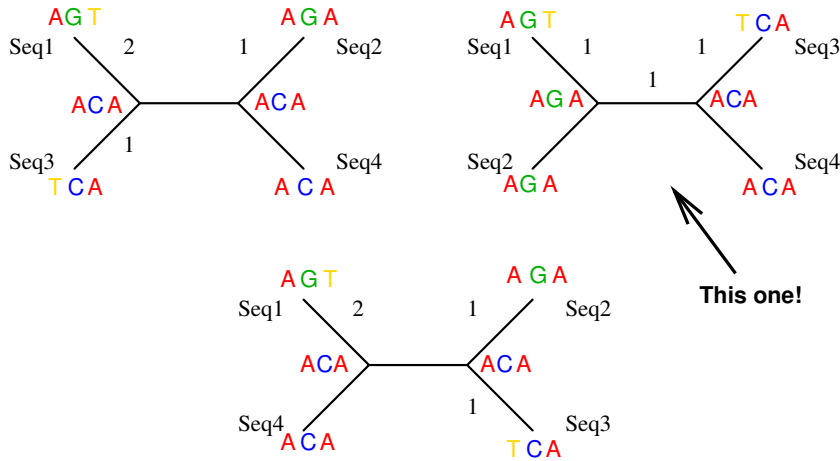
$x^1$	$=$	$x_1^1, x_2^1, \dots, x_m^1$	e.g. $n = 4$ :
			Seq1
			Seq2
			Seq3
			Seq4
$x^n$	$=$	$x_1^n, x_2^n, \dots, x_m^n$	

Which tree is *most parsimonious*, i.e. explains the data with the least number of mutations?

For this question we can neglect all non-polymorphic sites.

Which tree is most parsimonious?

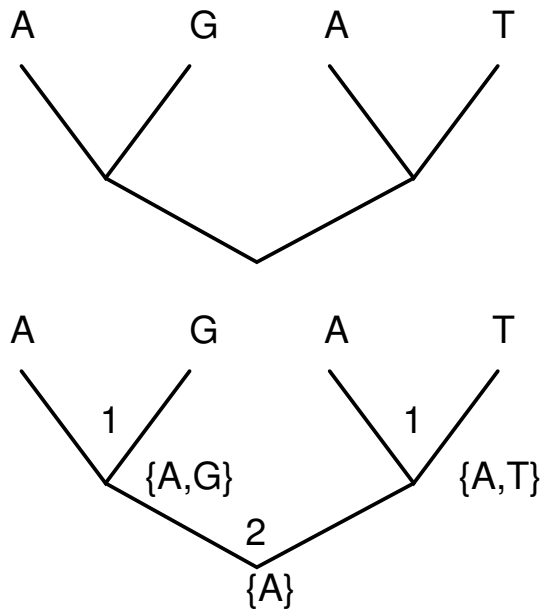




Given a tree whose tips are labeled with sequences, how can we efficiently compute the minimal number of mutations?

ideas:

1. Do separately for each alignment column
2. label each inner node with the optimal states for the tips above it and with the least number of mutations
3. go from tips to root by dynamic programming



### Fitch algorithm

$C$  is a counter of mutations, and  $M_k$  is the set of optimal states in node  $k$ .

Do for all sites  $s$ :

1.  $C_s := 0$  will be the counter of mutations at that site
2. for all tips  $b$  with label  $x$  set  $M_b = \{x\}$ .
3. Moving from tips to root do for all nodes  $k$  with daughter nodes  $i$  and  $j$ :
  - if**  $M_i \cap M_j = \emptyset$ : set  $M_k = M_i \cup M_j$  and  $C_s := C_s + 1$
  - else**: set  $M_k = M_i \cap M_j$

output  $\sum_s C_s$

### weighted parsimony

It is possible to take into account that different types of mutations (e.g. transitions and transversions) differ in the frequency by defining a cost  $S(a, b)$  for a mutation  $a \rightarrow b$ .

A variant of the Fitch algorithm calculates the minimal cost of a given tree to generate given sequences at the tips. ( $\rightsquigarrow$  exercise)

### 3.2 Finding parsimonious trees for given data

Given a large number  $n$  taxa, it is not feasible to consider all trees, because the number of unrooted bifurcating trees with  $n$  taxa is

$$3 \times 5 \times 7 \times \dots \times (2n - 5)$$

$n$	$3 \times 5 \times 7 \times \dots \times (2n - 5)$
5	15
7	945
10	2,027,025
12	654,729,075
20	$2.2 \cdot 10^{20}$
50	$2.8 \cdot 10^{74}$
100	$1.7 \cdot 10^{182}$

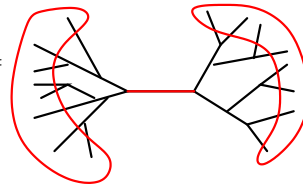
for comparison:

(estimated number of atoms in the observable universe)  $\times$  (number of second since big bang)  $\approx 5 \cdot 10^{97}$

#### problem of perfect phylogeny

Given  $n$  sequences of length  $m$  with up to 2 different states per position (alignment column). Is there a perfectly parsimonious tree, i.e. one that never has more than one mutation at the same position?

Idea: each polymorphism defines a split of the set of taxa  $L = A \cup B$ ,  $A \cap B = \emptyset$ .  
A branch of a tree also defines a split of  $L$



Go through the alignment from left to right and further subdivide  $L$  until there is a contradiction or you reach the end of the alignment.

**Theorem 3** A contradiction will occur if and only if there are two polymorphisms that lead to two splits  $L = A \cup B = C \cup D$  such that the four intersections  $A \cap C$ ,  $A \cap D$ ,  $B \cap C$ ,  $B \cap D$  are all non-empty.

This gives us an efficient solution for the problem of perfect parsimony. How about a slight generalization?

Given  $n$  homologous sequences of length  $m$  with up to  $r$  different states in each column.

Is there a perfectly parsimonious tree, i.e. one without back-mutations and without more than one mutation into the same state in the same position?

complexity: NP-complete for unbounded  $r$  and polynomial for any fixed  $r \in \mathbb{N}$ .

#### The problem of maximum parsimony

Given  $n$  homologous sequences of length  $m$  with up to 2 different states in each column, find the tree that needs the minimum number of mutations to explain the tree.

complexity: NP-complete

There is a method that can guarantee to find a tree that needs at most twice as many mutations as needed by the most parsimonious tree. However, in practice heuristic search algorithms are more relevant.

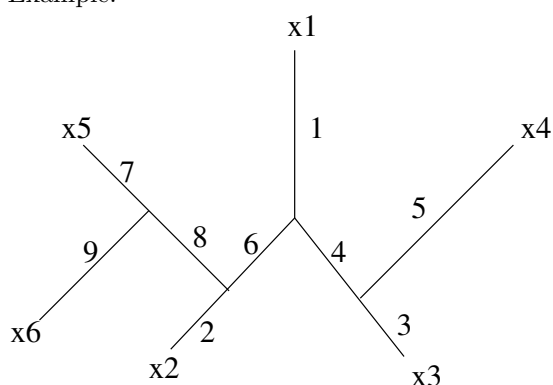
### enumerating all tree topologies

The sequence of numbers  $[i_3][i_5][i_7] \dots [i_{2n-5}]$  with  $i_k \in \{1, \dots, k\}$  represents a tree topology with  $n$  labeled leaves. It can be decoded as follows.

- Start with a 3-leaved tree whose leaves are labeled with  $x_1, x_2, x_3$  and whose edges are labeled accordingly with 1,2,3.
- repeat for  $j = 4, \dots, n$ :
  1.  $k := 2j - 5$
  2. Add an edge to the new leaf  $x_j$  to edge  $i_k$
  3. Call the new edge  $k + 2$ .
  4. In the subdivided edge  $i_k$ , give the part that is closer to  $x_1$  the label  $k + 1$ . The other part keeps the label  $i_k$ .

### enumerating all tree topologies

Example:



This tree can be represented by  $[3][2][7]$

### enumerating all labeled tree topologies

Enumerate leaves-labeled topologies by iterating  $[a][b][c] \dots [x]$  like a mileage counter for all allowed values ( $a \leq 3, b \leq 5, c \leq 7, \dots$ ):

$$\begin{array}{c}
 [1][1][1] \dots [1][1][1] \\
 [1][1][1] \dots [1][1][2] \\
 [1][1][1] \dots [1][1][3] \\
 \vdots \\
 [1][1][1] \dots [1][1][n] \\
 [1][1][1] \dots [1][2][1] \\
 [1][1][1] \dots [1][2][2] \\
 [1][1][1] \dots [1][2][3] \\
 \vdots
 \end{array}$$

### Branch and Bound

Let

$$[3][4][2] \dots [19][0][0][0]$$

denote the tree in which the last three taxa are not yet inserted. (zeros are only allowed at the end of a series).

Now we also iterate over these trees. If, e.g.  $u, v, w, x, y$  are the maxima of the last five positions:

$$\begin{array}{c}
 [a][b][c] \dots [m][u][v][w][x][y] \\
 [a][b][c] \dots [m][0][0][0][0][0] \\
 [a][b][c] \dots [m][1][0][0][0][0] \\
 [a][b][c] \dots [m][1][1][0][0][0] \\
 [a][b][c] \dots [m][1][1][1][0][0]
 \end{array}$$

If the tree corresponding to  $[a][b][c]...[m][1][1][1][0][0]$  already needs more mutations than the best tree found so far, go directly to

$$[a][b][c]...[m][1][1][2][0][0]$$

(“Bound”)

“Branch and Bound” saves time and can be used in practice for up to about 11 taxa.

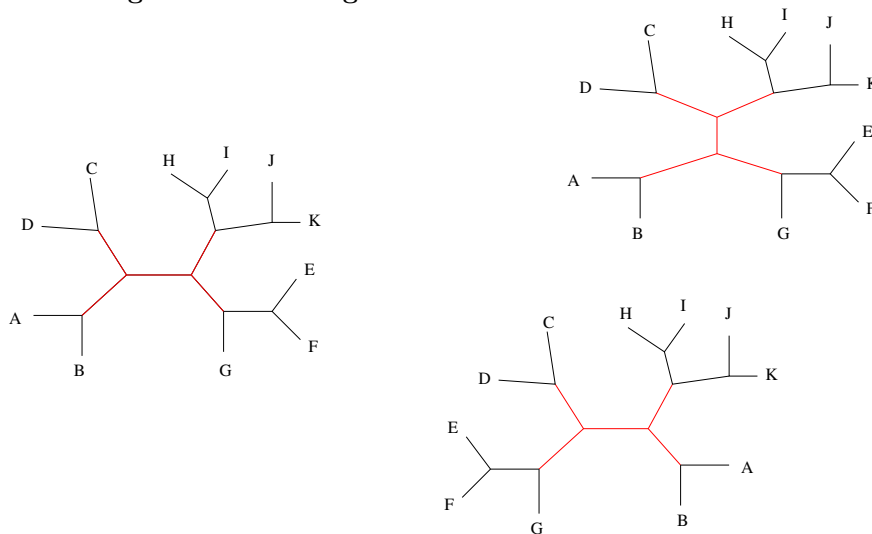
For larger numbers of taxa we need to move around in tree space and try to optimize the tree topology. Possible steps are

**NNI:** nearest neighbor interchange

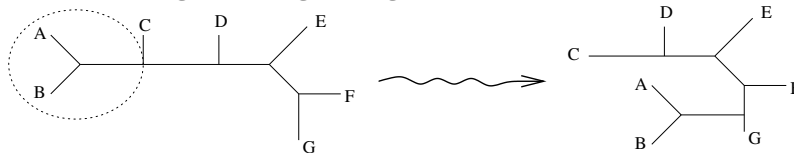
**SPR:** subtree pruning and regrafting

**TBR:** tree bisection and reconnection

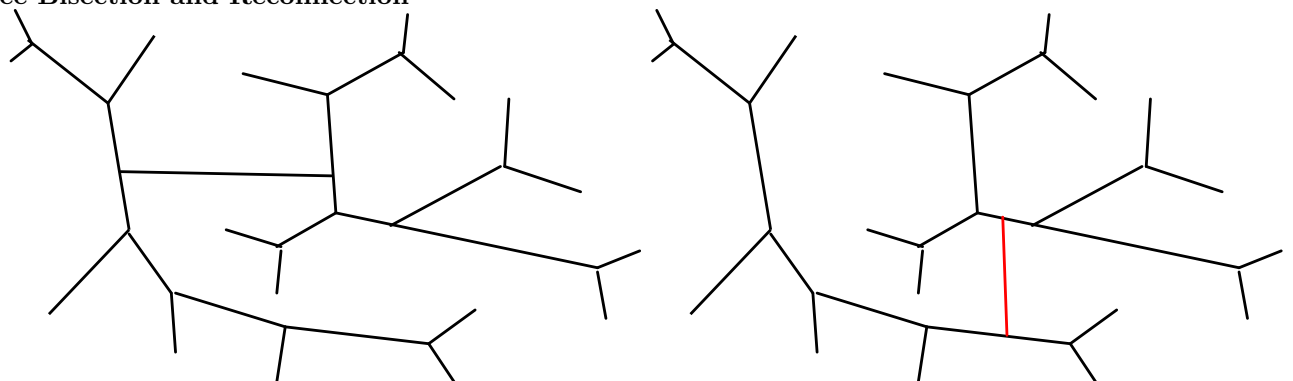
### Nearest Neighbor Interchange



### Subtree Pruning and Regrafting



### Tree Bisection and Reconnection



Note that each NNI move is special type of SPR move where the pruned subtree is regrafted in an edge neighboring the original edge.

Each SPR move is a special TBR where one of the nodes of the new edge is the old node.

### 3.3 Excursus: measuring the difference between two trees

#### The symmetric difference (aka “partition metric”)

Bourque (1978), Robinson and Foulds (1981)

Each edge in the tree is a partition of the set of taxa. The *symmetric difference* is the number of edges that exist in one tree but not in the other.

#### quartet distance

for fully resolved trees of  $n$  taxa.

Each of the  $\binom{n}{4}$  quartets of taxa have a tree topology in each tree. The *quartet distance* is the relative frequency of quartets for which the topologies do not coincide.

#### NNI distance

Waterman, Smith (1978)

The *NNI distance* is the number of NNI moves needed to change the one tree topology into the other.

Problem: It has been shown that the computation of the NNI distance is NP-hard.

Allen and Steel (2001) showed that the TBR distance is easier to compute.

#### Path-length difference metric

Penny, Watson, Steel 1993

Let  $n_{ab}^T$  be the number of edges that separate taxa  $a$  and  $b$  in tree  $T$ . Then, the path-length difference metric between the trees  $T$  and  $T'$  is defined as

$$\sqrt{\sum_{a,b} (n_{ab}^T - n_{ab}^{T'})^2}$$

#### taking branch lengths into account

For each partition  $P$  of the taxa set let  $f_T(P)$  be the length of the corresponding edge in tree  $T$  if such an edge exists. Otherwise set  $f_T(P) = 0$

**branch score distance** (Kuhner, Felsenstein, 1994)

$$\sum_P (f_T(P) - f_{T'}(P))^2$$

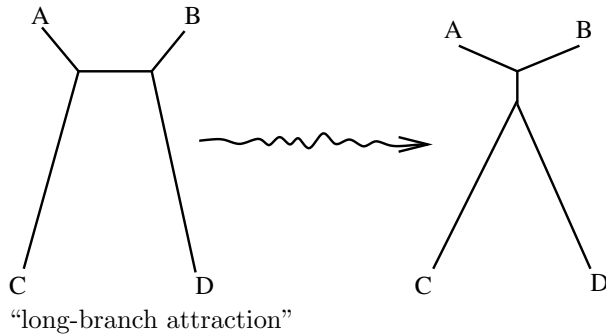
**Robinson-Foulds distance**

$$\sum_P |f_T(P) - f_{T'}(P)|$$

### 3.4 Limitations of the parsimony principle

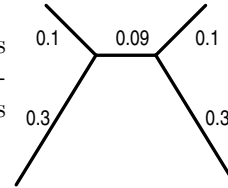
#### limitations of parsimony

Parsimonious phylogeny reconstruction methods do not take back mutations and double-hits into account in a proper way. This can lead to problems when there are long branches with many mutations.



### Comparison of phylogeny estimation methods

Durbin et al. (1998) simulated for several sequence lengths 1000 quartets of sequences along this tree to compare the accuracy of phylogeny reconstruction methods. Branch lengths are mean frequencies of transitions per position. (no transversions)



Proportion of correctly estimated trees			
Seq.length	Max.Pars.	Neigh.Join.	ML
20	39.6%	47.7%	41.9%
100	40.5%	63.5%	63.8%
500	40.4%	89.6%	90.4%
2000	35.3%	99.5%	99.7%

## 4 Maximum-Likelihood (ML) in phylogeny estimation

### 4.1 What is a likelihood?

#### Frequentistic parameter estimation

- Assume that we observe some data  $D$ .
- $D$  is influenced by random effects but a parameter  $p$  plays a role.
- We are interested in the value of  $p$
- $D$  is random but observed
- $p$  is unknown but not random
- A model describes how the probability of  $D$  depends on  $p$

Maximum-Likelihood principle: estimate  $p$  by

$$\hat{p} = \arg \max_p \Pr_p(D)$$

To describe how  $\Pr_p(D)$  depends on  $p$  we define the likelihood function:

$$L_D(p) := \Pr_p(D)$$

The ML estimator  $\hat{p}$  is the parameter value that maximizes the probability of the observed data.

### simple example

If you toss a thumbtack, what is the probability  $p$  that the sting touches the ground?  
 Assume you made an experiment. In 1000 tosses, the sting touched the ground 567 times.

$$\begin{aligned} \hat{p} &= \arg \max_p \Pr_p(567) = \arg \max_p \binom{1000}{567} p^{567} \cdot (1-p)^{1000-567} \\ &= \arg \max_p \log(p^{567} \cdot (1-p)^{433}) \\ &= \arg \max_p 567 \cdot \log(p) + 433 \cdot \log(1-p) \\ \\ \frac{\partial}{\partial p} (567 \log(p) + 433 \log(1-p)) &= \frac{567}{p} - \frac{433}{1-p} \\ \Rightarrow 0 = \frac{567}{\hat{p}} - \frac{433}{1-\hat{p}} &\Rightarrow \hat{p} = 0.567 \end{aligned}$$

**Important:** the parameter  $p$  is *not* a random object. Thus it does not make sense to ask for the *probability* that it takes some particular value  $p_0$ . However, the *likelihood* of  $p_0$  is defined. It is the probability of the observed data if  $p = p_0$ .

### 4.2 How to compute the likelihood of a tree

ML estimation of phylogenetic trees: Given an alignment  $D$ , find the tree  $T$  that maximizes

$$\Pr(D|T) =: L_D(T), \quad \hat{T} := \arg \max_T L_D(T)$$

What is  $\Pr(D|T)$  and how can we compute it?

We assume that all alignment columns evolve independently of each other. Then

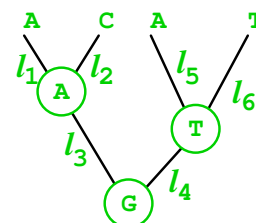
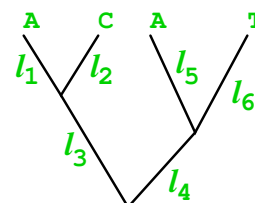
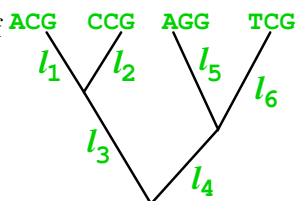
$$\Pr(D|T) = \prod_i \Pr(d_i|T),$$

where  $d_i$  is the sequence data in the  $i$ -th alignment column.

But how can we compute  $\Pr(d_i|T)$ ?

How to compute  $\Pr(d_i|T) = L_{d_i}(T)$ ?

Let's first assume that  $d_i$  also contains labels of the inner nodes. Assume that for all nucleotide  $x, y$  and all  $\ell \in \mathbb{R}_{>0}$  we can compute the frequency  $p_x$  of  $x$  and the probability  $P_{x \rightarrow y}(\ell)$  that an  $x$  is replaced by a  $y$  along a branch of length  $\ell$ .



Then, we get for the example tree

$$\begin{aligned} \Pr(d_i|T) &= p_G \cdot P_{G \rightarrow A}(\ell_3) \cdot P_{G \rightarrow T}(\ell_4) \cdot \\ &\quad \cdot P_{A \rightarrow A}(\ell_1) \cdot P_{A \rightarrow C}(\ell_2) \cdot \\ &\quad \cdot P_{T \rightarrow A}(\ell_5) \cdot P_{T \rightarrow T}(\ell_6). \end{aligned}$$

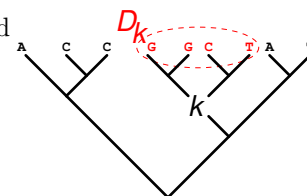
But usually, inner nodes are not labeled. What to do then?

### Felsenstein's pruning algorithm

For each node  $k$  let  $D_k$  be the part of the data  $d_i$  that are labeled to tips that stem from  $k$  and define

$$w_k(x) = \Pr(D_k | k \text{ has an } x \text{ at this site})$$

for every nucleotide  $x$ .



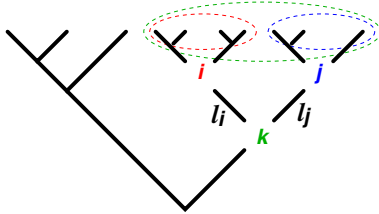
Idea: compute  $w_k(x)$  for all  $k$  and all  $x$ . Then you know it also for the root  $r$  and can compute

$$L(T) = \Pr(D|T) = \Pr(D_r|T) = \sum_{x \in \{A,C,G,T\}} p_x \cdot w_r(x).$$

Compute all  $w_k(x)$  from the tips to the root by dynamic programming.

For any leave  $b$  with nucleotide  $y$  we have

$$w_b(x) = \begin{cases} 0 & \text{if } x \neq y \\ 1 & \text{if } x = y \end{cases}$$



If  $k$  is a node with child nodes  $i$  and  $j$  and corresponding branch lengths  $\ell_i$  and  $\ell_j$ , then

$$w_k(x) = \left( \sum_{y \in \{A,C,G,T\}} P_{x \rightarrow y}(\ell_i) \cdot w_i(y) \right) \cdot \left( \sum_{z \in \{A,C,G,T\}} P_{x \rightarrow z}(\ell_j) \cdot w_j(z) \right)$$

How to compute  $P_{x \rightarrow y}(\ell)$ ?

You need a model for sequence evolution. The simplest one is the Jukes-Cantor model:

- all sites independent of each other (given the tree)
- all  $p_x$  equal
- “mutations” appear at rate  $\lambda$
- a “mutation” lets the site forget its state and sample the new one uniformly from  $\{A, C, G, T\}$ . (i.e.  $A$  can be replaced by another  $A$ )
- (in original paper for protein sequences)

### What is a rate?

Let  $M_{a,b}$  the number of “mutations” in time interval  $[a, b]$ .

- Rate  $\lambda$  means that the expected number of “mutations” in a time interval of length  $t$  is  $\lambda t$ :

$$\mathbb{E}M_{0,t} = \lambda t$$

- If  $\varepsilon > 0$  is extremely small, then we may neglect the probability of more than one “mutation” in a time interval of length  $\varepsilon$ .
- Then,  $\lambda \varepsilon$  is not only the expected number of mutations but also the probability that there is one in that time interval:

$$\Pr(M_{0,\varepsilon} > 0) \approx \Pr(M_{0,\varepsilon} = 1) \approx \mathbb{E}M_{0,\varepsilon} = \lambda \varepsilon$$

- numbers of “mutations” on disjoint intervals are stochastically independent

For longer time intervals  $[0, t]$  we choose a large  $n \in \mathbb{N}$  and argue:

$$\begin{aligned} \Pr(M_{0,t} = 0) &= \Pr(M_{0,t/n} = 0, M_{t/n, 2t/n} = 0, \dots, M_{(n-1)t/n, t} = 0) \\ &= \Pr(M_{0,t/n} = 0) \cdot \Pr(M_{t/n, 2t/n} = 0) \cdots \Pr(M_{(n-1)t/n, t} = 0) \\ &\approx \left(1 - \lambda \frac{t}{n}\right)^n \xrightarrow{n \rightarrow \infty} e^{-\lambda t} \end{aligned}$$



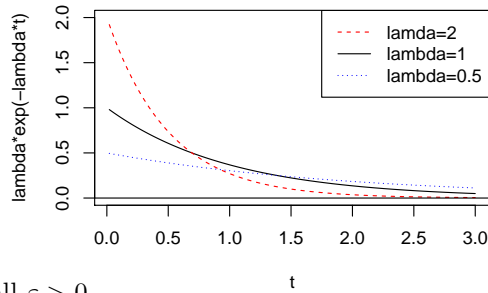
This means: the waiting time  $\tau$  for the first mutation is **exponentially distributed** with rate  $\lambda$ . This means it has

$$\Pr(\tau > t) = e^{-\lambda t}$$

**expectation value**  $\mathbb{E}\tau = 1/\lambda$

**standard deviation**  $\sigma_\tau = 1/\lambda$

**density**  $f(t) = \lambda \cdot e^{-\lambda t}$



This means  $\Pr(\tau \in [t, t + \epsilon]) \approx f(t) \cdot \epsilon$  for small  $\epsilon > 0$ .

After this preparation we can finally compute  $P_{x \rightarrow y}(t)$ , first for  $y \neq x$ :

$$\begin{aligned} P_{x \rightarrow y}(t) &= \Pr(M_{0,t} > 0) \cdot \Pr(\text{last "mutation" leads to } y) \\ &= (1 - e^{-\lambda t}) \cdot \frac{1}{4} \end{aligned}$$

and

$$\begin{aligned} P_{x \rightarrow x}(t) &= \Pr(M_{0,t} = 0) + \Pr(M_{0,t} > 0) \cdot \Pr(\text{last "mutation" leads to } x) \\ &= e^{-\lambda t} + (1 - e^{-\lambda t}) \cdot \frac{1}{4} \\ &= \frac{1}{4} + \frac{3}{4}e^{-\lambda t} \end{aligned}$$

### 4.3 How to search for the ML tree

Given a large number  $n$  of taxa (i.e. sequences), it is difficult to find the ML phylogeny. Two partial problems have to be solved:

1. Given the tree topology, find the optimal branch lengths
2. Find the tree topology for which your solution of problem 1 leads to the highest likelihood value.

We first turn to problem 1.

#### Tree length optimization in the very first version of PHYLIP dnaml

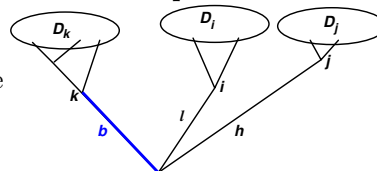
Expectation-Maximization (EM) algorithm: Iterate the following steps:

**E step** given the current branch lengths and rates, compute the expected number of mutations for each branch

**M step** optimize branch lengths for the expected numbers of mutations computed in the E step

**More common: use the derivative of the likelihood with respect to the branch length**

To optimize the length  $b$  of some branch, first rotate it, such that one of its adjacent nodes is the root.



First we assume that the alignment  $D$  has only one column. Then,  $\hat{L}_D(T)$  is  $\sum_x p_x \cdot \sum_y P_{x \rightarrow y}(b) \cdot w_k(y) \cdot (\sum_z P_{x \rightarrow z} w_i(z)) \cdot (\sum_{z'} P_{x \rightarrow z'} w_j(z'))$ .

$$\Rightarrow \frac{\partial L_D(T)}{\partial b} = \sum_x p_x \cdot \sum_y \frac{\partial P_{x \rightarrow y}(b)}{\partial b} \cdot w_k(y) \cdot (\sum_z P_{x \rightarrow z} w_i(z)) \cdot (\sum_{z'} P_{x \rightarrow z'} w_j(z'))$$

and  $\frac{\partial^2 L_D(T)}{\partial b^2} = \sum_x p_x \cdot \sum_y \frac{\partial^2 P_{x \rightarrow y}(b)}{\partial b^2} \cdot w_k(y) \cdot (\sum_z P_{x \rightarrow z} w_i(z)) \cdot (\sum_{z'} P_{x \rightarrow z'} w_j(z'))$   
 In the Jukes-Cantor model we can compute for example for  $x \neq y$ :

$$\frac{\partial}{\partial b} P_{x \rightarrow y}(b) = \frac{\partial}{\partial b} (1 - e^{-\lambda b}) \cdot \frac{1}{4} = \frac{1}{4} \lambda e^{-\lambda b}$$

$$\frac{\partial^2}{\partial b^2} P_{x \rightarrow y}(b) = -\frac{1}{4} \lambda^2 e^{-\lambda b}$$

For alignments  $D$  with columns  $D_1 \dots D_m$  we can compute all  $L'_h := \frac{\partial}{\partial b} L_{D_h}(T)$  and  $L''_h := \frac{\partial^2}{\partial b^2} L_{D_h}(T)$  as explained above, and then compute the first two derivatives of  $L_D(T) = \prod_h L_{D_h}(T)$  by applying the product rule for derivatives:

$$\frac{\partial}{\partial b} L_D(T) = L_D(T) \cdot \sum_h \frac{L'_h}{L_{D_h}(T)}$$

and

$$\frac{\partial^2}{\partial b^2} L_D(T) = L_D(T) \cdot \sum_h \left( \frac{L''_h}{L_{D_h}(T)} + \sum_{\ell \neq h} \frac{L'_h \cdot L'_\ell}{L_{D_h}(T) \cdot L_{D_\ell}(T)} \right)$$

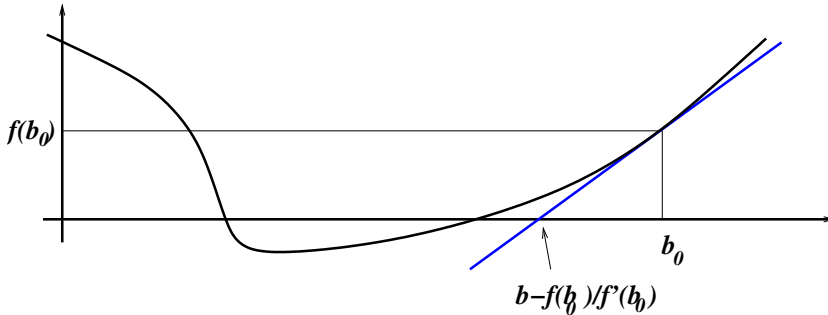
To optimize  $b$ , solve

$$f(b) := \frac{\partial L_D(T)}{\partial b} = 0.$$

This is numerically with a Newton-Raphson scheme using  $f'(b) = \frac{\partial^2 L_D(T)}{\partial b^2}$ .

### Newton-Raphson scheme to solve $f(b) = 0$

1. Start with some initial value  $b_0$
2. as long as  $f(b_0)$  is not close enough to 0, replace  $b_0$  by  $b_0 - f(b_0)/f'(b_0)$  and try again.



### Optimizing the topology

Now that know how to search for the optimal tree, given the topology, how do we search for the best topology?

[stepwise addition](#) (default of DNAML):

- start with the only possible tree of the first three taxa
- stepwise add one taxon
- to do this when  $k$  taxa are already added, try all  $2k - 5$  possible branches to add the next taxon, optimize branch lengths
- when all are added, optimize with NNI steps
- repeat whole procedure with different input orders

[branch and bound](#) if only few taxa

[Start with NeighborJoining and continue with SPR](#) is nowadays most common

[Supertree methods](#) like PUZZLE: ML for all quartets, then build tree that respects most of them.

## 4.4 Consistency of the Maximum-Likelihood method

**Theorem 4** *The ML estimator for phylogenetic trees is **consistent**. This means, if the model assumptions are fulfilled and you add more and more data (i.e. make the sequences longer) for a fixed set of taxa, the probability that the ML tree will converge against the true tree is 1.*

Notes:

1. the ML tree is the tree with the highest likelihood. ML tree estimation programs do not always find the ML tree
2. the model assumptions include a model for the substitution process and that all sequence positions are independent and correctly aligned

Sketch of proof for the consistency of the ML tree:

Let  $a_1, \dots, a_m$  be the different alignment columns and let  $n_1, \dots, n_m$  be their numbers of occurrences in the data  $D$ . The likelihood of a tree  $T$  is then

$$L(T) = \Pr(D | T) = \prod_{i=1}^m \Pr(a_i | T)^{n_i}$$

Idea of the proof: the probabilities  $\Pr(a_i | T)$  are characteristic for  $T$ , and those of the true tree will be reflected in the relative frequencies  $R_i = n_i/n$  with  $n = n_1 + \dots + n_m$ .

The log likelihood is

$$\ln L(T) = \sum_{i=1}^m n_i \ln \Pr(a_i | T) = n \cdot \sum_{i=1}^m R_i \ln \Pr(a_i | T)$$

For long sequences we get  $R_i \rightarrow p_i$ , where  $p_1, \dots, p_m$  are the (unknown) probabilities of  $a_1, \dots, a_m$  for the true tree  $T^*$ . Let  $q_1, \dots, q_m$  be those probabilities for some other tree  $T$ . Then we obtain

$$\frac{1}{n} \ln L(T) \xrightarrow{n \rightarrow \infty} \sum_{i=1}^m p_i \ln q_i \quad \text{and} \quad \frac{1}{n} \ln L(T^*) \xrightarrow{n \rightarrow \infty} \sum_{i=1}^m p_i \ln p_i.$$

For  $p \neq q$  we get

$$\sum_{i=1}^m p_i \ln q_i < \sum_{i=1}^m p_i \ln p_i,$$

because

$$\sum_{i=1}^m p_i \ln p_i - \sum_{i=1}^m p_i \ln q_i = \sum_{i=1}^m p_i \ln \frac{p_i}{q_i} > 0,$$

and the last inequation follows since  $\sum_{i=1}^m p_i \ln \frac{p_i}{q_i}$  is the relative entropy, also called Kullback-Leibler-Information, which is known to be positive for  $p \neq q$ . (Taking partial derivatives  $\frac{\partial}{\partial q_i}$  with boundary condition  $\sum q_i = 1$  shows that  $\sum_{i=1}^m p_i \ln \frac{p_i}{q_i} = 0$  is the minimum.)

## 4.5 Maximum Parsimony from a probabilistic perspective

If we assume a probabilistic substitution model, we can set  $s(a, b) = -\log P_{a \rightarrow b}(1)$  and use the values  $s(a, b)$  as costs in weighted parsimony. Thus, maximum parsimony can be considered as an approximation for the case that

1. all edges in the tree have the same length
2. double-hits and back-mutations are negligible

## 4.6 Maximum likelihood for pairwise distances

Given a substitution model with known parameters we can compute the ML distance  $d_{xy}^{\text{ML}}$  between sequence  $x = (x_1, x_2, \dots, x_n)$  and sequence  $y = (y_1, \dots, y_n)$  by

$$\begin{aligned} d_{xy}^{\text{ML}} &= \arg \max_t \left\{ \prod_i \pi_{x_i} \cdot P_{x_i \rightarrow y_i}(t) \right\} \\ &= \arg \max_t \left\{ \prod_i P_{x_i \rightarrow y_i}(t) \right\} \end{aligned}$$

E.g. for the Jukes-Cantor Model with rate  $\alpha$  we get in the case of  $k$  mismatches:

$$\prod_i P_{x_i \rightarrow y_i}(t) = \left( \frac{1}{4}(1 - e^{-t\alpha}) \right)^k \left( \frac{1}{4}(1 + 3e^{-t\alpha}) \right)^{n-k}$$

Optimizing this with the usual procedure we get:

$$d_{xy}^{\text{ML}} = -\frac{1}{4\alpha} \cdot \ln \left( 1 - \frac{4k}{3n} \right)$$

Also this ML estimator is consistent, i.e. will give us the true distances in the limit of long sequences. This implies that applying NeighborJoining to the ML distances is also consistent.

If the sequences are not extremely long, direct ML methods will tend to give more reliable results (as long as they are computationally tractable.)

## 5 Modelling the substitution process on sequences

### 5.1 Transition matrix and rate matrix

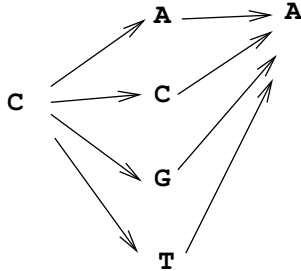
Let  $P_{a \rightarrow b}(t)$  be the probability that a nucleotide  $a$  is a nucleotide  $b$  after time (i.e. branch length)  $t$ .

$$S(t) := \begin{pmatrix} P_{A \rightarrow A}(t) & P_{A \rightarrow C}(t) & P_{A \rightarrow G}(t) & P_{A \rightarrow T}(t) \\ P_{C \rightarrow A}(t) & P_{C \rightarrow C}(t) & P_{C \rightarrow G}(t) & P_{C \rightarrow T}(t) \\ P_{G \rightarrow A}(t) & P_{G \rightarrow C}(t) & P_{G \rightarrow G}(t) & P_{G \rightarrow T}(t) \\ P_{T \rightarrow A}(t) & P_{T \rightarrow C}(t) & P_{T \rightarrow G}(t) & P_{T \rightarrow T}(t) \end{pmatrix}$$

Each row has sum 1.

How can we compute  $S(2)$  from  $S(1)$ ?

For example:  $P_{C \rightarrow A}(2)$



$$\begin{aligned} P_{C \rightarrow A}(2) &= P_{C \rightarrow A}(1) \cdot P_{A \rightarrow A}(1) + \\ &P_{C \rightarrow C}(1) \cdot P_{C \rightarrow A}(1) + \\ &P_{C \rightarrow G}(1) \cdot P_{G \rightarrow A}(1) + \\ &P_{C \rightarrow T}(1) \cdot P_{T \rightarrow A}(1) \end{aligned}$$

With matrix multiplication we can write this as

$$P(2) = P(1) \cdot P(1).$$

More generally:

$$P(t + s) = P(t) \cdot P(s)$$

**Matrix multiplication**  $A \cdot B = C$

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \ddots & & \vdots \\ a_{i1} & a_{i2} & \cdots & a_{im} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{pmatrix} \cdot \begin{pmatrix} b_{11} & b_{12} & \cdots & b_{1j} & \cdots & b_{1k} \\ b_{21} & b_{22} & \cdots & b_{2j} & \cdots & b_{2k} \\ \vdots & \ddots & & \vdots & \ddots & \vdots \\ b_{m1} & b_{m2} & \cdots & b_{mj} & \cdots & b_{mk} \end{pmatrix} \\ = \begin{pmatrix} c_{11} & c_{12} & \cdots & \cdots & c_{1k} \\ c_{21} & c_{22} & \cdots & \cdots & c_{2k} \\ \vdots & \ddots & & c_{ij} & \vdots \\ c_{n1} & c_{n2} & \cdots & \cdots & c_{nk} \end{pmatrix}, \quad c_{ij} = \sum_{h=1}^m a_{ih} \cdot b_{hj}$$

**Matrix addition**  $A + B = C$

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{pmatrix} + \begin{pmatrix} b_{11} & b_{12} & \cdots & b_{1m} \\ b_{21} & b_{22} & \cdots & b_{2m} \\ \vdots & \ddots & & \vdots \\ b_{n1} & b_{n2} & \cdots & b_{nm} \end{pmatrix} \\ = \begin{pmatrix} a_{11} + b_{11} & a_{12} + b_{12} & \cdots & a_{1m} + b_{1m} \\ a_{21} + b_{21} & a_{22} + b_{22} & \cdots & a_{2m} + b_{2m} \\ \vdots & \ddots & & \vdots \\ a_{n1} + b_{n1} & a_{n2} + b_{n2} & \cdots & a_{nm} + b_{nm} \end{pmatrix}$$

Rules:

$$A + B = B + A \quad (A + B) + C = A + (B + C) \quad (A \cdot B) \cdot C = A \cdot (B \cdot C) \\ A \cdot (B + C) = (A \cdot B) + (A \cdot C) \quad \text{but in general } A \cdot B \neq B \cdot A$$

For very small  $\varepsilon > 0$ ,  $P_{x \rightarrow x}(\varepsilon)$  is close to 1 and there is a matrix  $R$ , the so-called **rate matrix** (or  $Q$ -matrix), such that  $S(\varepsilon) \approx (I + R \cdot \varepsilon)$ , where  $I$  is the identity matrix (or unit matrix)

$$I = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1 \end{pmatrix}$$

with the property that  $A \cdot I = A$  and  $I \cdot B = B$  for all matrices  $A$  and  $B$  of suitable dimensions.

Thus, we obtain  $S(t + \varepsilon) = S(t) \cdot S(\varepsilon) \approx S(t)(I + R\varepsilon) = S(t) + S(t)R\varepsilon$  and

$$\lim_{\varepsilon \rightarrow 0} \frac{S(t + \varepsilon) - S(t)}{\varepsilon} = S(t)R$$

$S(t)R$  is like the derivative of the process. Note that the row sums are 0. The diagonal entries are negative. All other entries are the rates of the corresponding substitutions.

Rate matrix of the Jukes-Cantor-Model for DNA

$$\begin{pmatrix} -\frac{3}{4}\alpha & \frac{1}{4}\alpha & \frac{1}{4}\alpha & \frac{1}{4}\alpha \\ \frac{1}{4}\alpha & -\frac{3}{4}\alpha & \frac{1}{4}\alpha & \frac{1}{4}\alpha \\ \frac{1}{4}\alpha & \frac{1}{4}\alpha & -\frac{3}{4}\alpha & \frac{1}{4}\alpha \\ \frac{1}{4}\alpha & \frac{1}{4}\alpha & \frac{1}{4}\alpha & -\frac{3}{4}\alpha \end{pmatrix}.$$

The model F81 (Felsenstein, 1981) allows for unequal nucleotide frequencies  $(\pi_A, \pi_C, \pi_G, \pi_T)$  and has the rate matrix

$$\begin{pmatrix} -\alpha + \alpha\pi_A & \alpha\pi_C & \alpha\pi_G & \alpha\pi_T \\ \alpha\pi_A & -\alpha + \alpha\pi_C & \alpha\pi_G & \alpha\pi_T \\ \alpha\pi_A & \alpha\pi_C & -\alpha + \alpha\pi_G & \alpha\pi_T \\ \alpha\pi_A & \alpha\pi_C & \alpha\pi_G & -\alpha + \alpha\pi_T \end{pmatrix}.$$

In addition, the HKY model (Hasegawa, Kishino, Yano, 1985) allows that transitions are more probable than transversions by using an additional parameter  $\beta$ . Its rate matrix is

$$Q := \begin{pmatrix} -\alpha\pi_G - \beta(\pi_C + \pi_T) & \beta\pi_A & \alpha\pi_G & \beta\pi_T \\ \beta\pi_A & -\alpha\pi_T - \beta(\pi_A + \pi_G) & \beta\pi_C & \alpha\pi_T \\ \alpha\pi_A & \beta\pi_C & -\alpha\pi_A - \beta(\pi_C + \pi_T) & \beta\pi_T \\ \beta\pi_A & \alpha\pi_C & \beta\pi_G & -\alpha\pi_C - \beta(\pi_A + \pi_G) \end{pmatrix}.$$

$(\pi_A, \pi_C, \pi_G, \pi_T)$  is the **stationary Distribution** (or **equilibrium distribution**) for any of these rate matrices. This means:

$$(\pi_A, \pi_C, \pi_G, \pi_T) \cdot Q = (0, 0, 0, 0)$$

Thus, if  $(\pi_A, \pi_C, \pi_G, \pi_T)$  are the nucleotide frequencies of the ancestral sequence, these frequencies are maintained if the sequence evolves according to HKY:

$$(\pi_A, \pi_C, \pi_G, \pi_T) \cdot S(t) = (\pi_A, \pi_C, \pi_G, \pi_T)$$

This also holds for the F81 model.

## 5.2 Residence time

If we think of discrete generations and a per generation mutation probability of  $p$ , the probability of seeing the first mutation in generation  $k$  is  $(1-p)^{k-1} \cdot p$ .

A random variable  $X$  with values in  $\{1, 2, \dots\}$  is **geometrically distributed** if  $\Pr(X = k) = (1-p)^{k-1} \cdot p$ .

Then,

$$\mathbb{E}X = \sum_{k=1}^{\infty} k \cdot (1-p)^{k-1} \cdot p = \frac{1}{p}$$

It is easy to check that this is the only possible value:

$$\begin{aligned} \mathbb{E}X &= \sum_{k=0}^{\infty} (k+1) \cdot (1-p)^k \cdot p \\ &= \sum_{k=1}^{\infty} k \cdot (1-p)^k \cdot p + \sum_{k=0}^{\infty} (1-p)^k \cdot p = (1-p) \cdot \mathbb{E}X + p \cdot \frac{1}{p} \\ \Rightarrow \quad \mathbb{E}X &= \frac{1}{p} \end{aligned}$$

The geometric distribution is characterized by the no-memory condition:

$$\Pr(X = k+n \mid X > k) = \Pr(X = n)$$

The continuous analogon is the exponential distribution: A random variable  $Y$  with values in  $\mathbb{R}_{\geq 0}$  is exponentially distributed with rate  $\lambda$  if

$$\Pr(Y > z) = e^{-\lambda z}.$$

In this case

$$\mathbb{E}Y = \int_0^{\infty} z \lambda e^{-\lambda z} = \frac{1}{\lambda}.$$

The exponential distribution approximates the geometric distribution if  $p$  is small and  $k$  is large:

$$(1 - p)^k \approx e^{-pk}.$$

In a continuous-time substitution model, the residence time in a state is exponential. For example, if a site has nucleotide **A**, and the HKY model applies, it stays an **A** for an exponentially distributed time with expectation value  $1/(\alpha\pi_G + \beta(\pi_C + \pi_T))$ . When it then mutates, it becomes a

$$\begin{aligned} \text{C} & \text{ with prob. } \frac{\beta\pi_C}{\alpha\pi_G + \beta(\pi_C + \pi_T)} \\ \text{G} & \text{ with prob. } \frac{\alpha\pi_G}{\alpha\pi_G + \beta(\pi_C + \pi_T)} \\ \text{T} & \text{ with prob. } \frac{\beta\pi_T}{\alpha\pi_G + \beta(\pi_C + \pi_T)}. \end{aligned}$$

### 5.3 Computing $S(t)$ from the rate matrix $R$

If  $S(1)$  is known, you can compute  $S(n)$  by

$$S(n) = S(1)^n.$$

To do this efficiently, diagonalize  $S(1)$ . This means, find a matrix  $U$  and a diagonal matrix  $D$  (this means  $D_{ij} = 0$  if  $i \neq j$ ), such that

$$S(1) = U \cdot D \cdot U^{-1}.$$

The inverse  $U^{-1}$  of the matrix  $U$  is defined by  $U^{-1} \cdot U = I = U \cdot U^{-1}$ .

In this case, the diagonal entries  $D_{ii} = \lambda_i$  of  $D$  are the eigenvalues of  $S(1)$ , the columns of  $U$  are the corresponding right eigenvectors and the rows of  $U^{-1}$  (with entries  $u'_{ij}$ ) are the left eigenvectors:

$$S(1) \cdot \begin{pmatrix} u_{1i} \\ u_{2i} \\ u_{3i} \\ u_{4i} \end{pmatrix} = \begin{pmatrix} u_{1i} \\ u_{2i} \\ u_{3i} \\ u_{4i} \end{pmatrix} \cdot \lambda_i, \quad (u'_{i1}, \dots, u'_{i4}) \cdot S(1) = \lambda_i \cdot (u'_{i1}, \dots, u'_{i4})$$

For

$$D = \begin{pmatrix} \mu_1 & 0 & \dots & 0 \\ 0 & \mu_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \mu_m \end{pmatrix}$$

we can use

$$D^n = \begin{pmatrix} \mu_1^n & 0 & \dots & 0 \\ 0 & \mu_2^n & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \mu_m^n \end{pmatrix}$$

and

$$\begin{aligned} S(1)^n &= (U \cdot D \cdot U^{-1})^n \\ &= U \cdot D \cdot U^{-1} \cdot U \cdot D \cdot U^{-1} \dots U \cdot D \cdot U^{-1} \cdot U \cdot D \cdot U^{-1} \\ &= U D \cdot I \cdot D \cdot I \dots D \cdot U^{-1} = U D^n U^{-1} \end{aligned}$$

The situation is similar in the continuous case. For  $t \in [0, \infty)$  we get  $S(t) = U \cdot T^t \cdot U^{-1}$  with

$$T^t = \begin{pmatrix} e^{\lambda_1 t} & 0 & \dots & 0 \\ 0 & e^{\lambda_2 t} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & e^{\lambda_m t} \end{pmatrix},$$

where  $\lambda_1, \lambda_2, \dots, \lambda_m$  are the eigenvalues of  $R$  (with  $m = 4$  for nucleotides and  $m = 20$  for amino acids).

Explanation: For very small  $\varepsilon > 0$  we have

$$S(t) = S(\varepsilon)^{t/\varepsilon} \approx (I + R \cdot \varepsilon)^{t/\varepsilon} = U_\varepsilon \cdot D_\varepsilon^{t/\varepsilon} \cdot U_\varepsilon^{-1},$$

where  $D_\varepsilon$  is a diagonal matrix of the eigenvalues  $\mu_i$  of  $I + \varepsilon \cdot R$ .

It is common to write this as  $S(t) = e^{tR}$  and call it “Matrix exponential”.

For the right eigenvectors  $v_i$  we have

$$(I + \varepsilon \cdot R) \cdot v_i = \mu_i \cdot v_i$$

and thus

$$R \cdot v_i = \frac{\mu_i - 1}{\varepsilon} \cdot v_i.$$

Therefore,

$$\lambda_i := \frac{\mu_i - 1}{\varepsilon}$$

is an eigenvalue of  $R$  (if  $\mu_i \neq 1$ ) and we can write the diagonal entries of  $D_\varepsilon^{t/\varepsilon}$  as

$$(\varepsilon\lambda_i + 1)^{t/\varepsilon},$$

which converges to  $e^{\lambda_i t}$  for  $\varepsilon \rightarrow 0$ .

Also check that rows of  $U$  and columns of  $U^{-1}$  are eigenvectors of  $I + R$  and thus also of  $R$ .

Since most programming languages provide functions for computing eigenvalues and eigenvectors, we can compute matrix exponentials. However, this is sometimes numerically instable.

One alternative is to use the following alternative definition of the matrix exponential:

$$e^{tR} = \sum_{n=0}^{\infty} \frac{(tR)^n}{n!}$$

which can be made more stable by choosing  $\beta > \max\{\lambda_1, \dots, \lambda_m\}$  and then using the variant

$$e^{tR} = e^{-\beta t} \cdot \sum_{n=0}^{\infty} \frac{(\beta t)^n \cdot (I + R/\beta)^n}{n!}.$$

Another approach is to use the limit

$$e^{tR} = \lim_{n \rightarrow \infty} \left( I + \frac{t}{n} R \right)^n$$

or its variant

$$e^{tR} = \lim_{n \rightarrow \infty} \left( \left( I - \frac{t}{n} R \right)^{-1} \right)^n$$

for the approximation

$$e^{tR} \approx \left( I + \frac{t}{n} R \right)^n$$

or

$$e^{tR} \approx \left( \left( I - \frac{t}{n} R \right)^{-1} \right)^n$$

with a large value of  $n$ .



## 5.4 A model for transition probabilities in closed form

The F84 model (Felsenstein, 1984) is similar to the HKY model but allows the computation of transition probabilities without numerics by using similar ideas as in the Jukes-Cantor model.

F84 model: Pepper crosses and bullets into the ancestral lineages of the all positions that make them (partly) forget their former type.

**crosses** come rate  $\lambda$ . The new type is drawn according to  $(\pi_A, \pi_C, \pi_G, \pi_T)$ .

**bullets** come at rate  $\mu$ . The lineage only remembers if it was a purine or a pyrimidine. If it was a purine, the new type is A or G with probability  $\frac{\pi_A}{\pi_A + \pi_G}$  or  $\frac{\pi_G}{\pi_A + \pi_G}$ . If it was a pyrimidine, the new type is C or T with probability  $\frac{\pi_C}{\pi_C + \pi_T}$  or  $\frac{\pi_T}{\pi_C + \pi_T}$ .

A transversion needs at least one cross. If we condition on having at least one cross but not on the nucleotide that was selected at the cross, then the last bullet or cross before time  $t$  draws a nucleotide according to the distribution  $(\pi_A, \pi_C, \pi_G, \pi_T)$ . Thus, we get, for example:

$$P_{A \rightarrow C}(t) = (1 - e^{-\lambda t}) \cdot \pi_C$$

A transition needs either at least one cross or no cross and at least one bullet. We get, for example:

$$P_{A \rightarrow G}(t) = (1 - e^{-\lambda t}) \cdot \pi_G + e^{-\lambda t} (1 - e^{-\mu t}) \cdot \pi_G / (\pi_A + \pi_G)$$

Even if we do not need it for computing the transition probabilities, we can write down the F84 rate matrix:

$$\begin{pmatrix} -\lambda(1 - \pi_A) - \frac{\mu\pi_G}{\pi_A + \pi_G} & \lambda\pi_C & \lambda\pi_G + \frac{\mu\pi_G}{\pi_A + \pi_G} & \lambda\pi_T \\ \lambda\pi_A & -\lambda(1 - \pi_C) - \frac{\mu\pi_T}{\pi_C + \pi_T} & \lambda\pi_G & \lambda\pi_T + \frac{\mu\pi_T}{\pi_C + \pi_T} \\ \lambda\pi_A + \frac{\mu\pi_A}{\pi_A + \pi_G} & \lambda\pi_C & -\lambda(1 - \pi_G) - \frac{\mu\pi_A}{\pi_A + \pi_G} & \lambda\pi_T \\ \lambda\pi_A & \lambda\pi_C + \frac{\mu\pi_C}{\pi_C + \pi_T} & \lambda\pi_G & -\lambda(1 - \pi_T) - \frac{\mu\pi_C}{\pi_C + \pi_T} \end{pmatrix}$$

## 5.5 Overview of DNA substitution models

**Jukes-Cantor-Modell (JC):** nucleotide type not considered.

from \ to	A	C	G	T
A	—	$\alpha$	$\alpha$	$\alpha$
C	$\alpha$	—	$\alpha$	$\alpha$
G	$\alpha$	$\alpha$	—	$\alpha$
T	$\alpha$	$\alpha$	$\alpha$	—

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

**Kimura's 2 Parameter Model (K2)** transitions more frequent than transversions.

from \ to	A	C	G	T
A	—	$\alpha$	$\beta$	$\alpha$
C	$\alpha$	—	$\alpha$	$\beta$
G	$\beta$	$\alpha$	—	$\alpha$
T	$\alpha$	$\beta$	$\alpha$	—

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

Felsenstein (1981) (**F81**) takes nucleotide frequencies  $(\pi_A, \pi_C, \pi_G, \pi_T)$  into account.

from \ to	A	C	G	T
A	—	$\alpha\pi_C$	$\alpha\pi_G$	$\alpha\pi_T$
C	$\alpha\pi_A$	—	$\alpha\pi_G$	$\alpha\pi_T$
G	$\alpha\pi_A$	$\alpha\pi_C$	—	$\alpha\pi_T$
T	$\alpha\pi_A$	$\alpha\pi_C$	$\alpha\pi_G$	—

**Hasegawa, Kishino und Yano (HKY)** regards nucleotide frequencies as well as differences between transitions and transversions.

from \ to	A	C	G	T
A	—	$\alpha\pi_C$	$\beta\pi_G$	$\alpha\pi_T$
C	$\alpha\pi_A$	—	$\alpha\pi_G$	$\beta\pi_T$
G	$\beta\pi_A$	$\alpha\pi_C$	—	$\alpha\pi_T$
T	$\alpha\pi_A$	$\beta\pi_C$	$\alpha\pi_G$	—

Felsenstein (1984) (**F84**) also regards nucleotide frequencies and differences between transitions and transversions. No matrix algebra is needed to compute transition probabilities,

from \ to	A	C	G	T
A	—	$\lambda\pi_C$	$\lambda\pi_G + \frac{\mu\pi_G}{\pi_A + \pi_G}$	$\lambda\pi_T$
C	$\lambda\pi_A$	—	$\lambda\pi_G$	$\lambda\pi_T + \frac{\mu\pi_T}{\pi_C + \pi_T}$
G	$\lambda\pi_A + \frac{\mu\pi_A}{\pi_A + \pi_G}$	$\lambda\pi_C$	—	$\lambda\pi_T$
T	$\lambda\pi_A$	$\lambda\pi_C + \frac{\mu\pi_C}{\pi_C + \pi_T}$	$\lambda\pi_G$	—

The **General Time-Reversible Model (GTR)** considers differences between pairs of nucleotide types.

von \ nach	A	C	G	T
A	—	$\alpha\pi_C$	$\beta\pi_G$	$\gamma\pi_T$
C	$\alpha\pi_A$	—	$\delta\pi_G$	$\epsilon\pi_T$
G	$\beta\pi_A$	$\delta\pi_C$	—	$\eta\pi_T$
T	$\gamma\pi_A$	$\epsilon\pi_C$	$\eta\pi_G$	—

In the models F81, F84, HKY and GTR,  $(\pi_A, \pi_C, \pi_G, \pi_T)$  is the stationary distribution, in JC and K2  $(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$ . All these models are reversible.

## 5.6 Position-dependent mutation rates

### Model for site-dependent rates

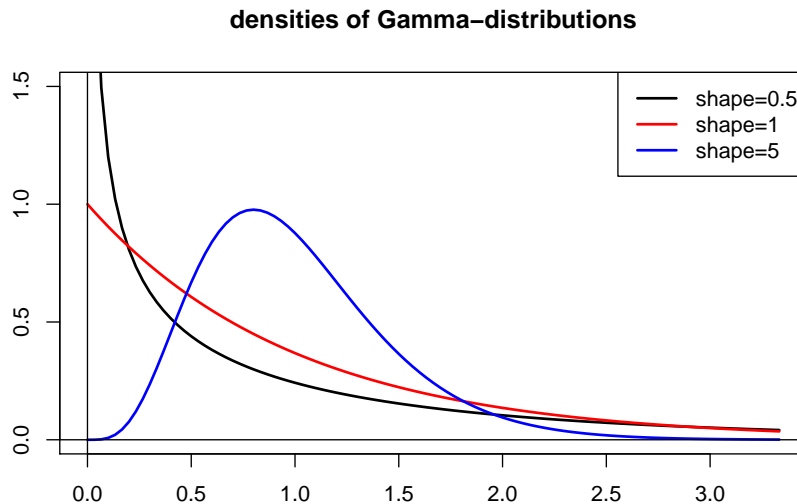
There is one rate matrix  $Q$  and for each site  $i$  there is a coefficient  $r_i$ , such that

$$R_i = r_i \cdot Q$$

is the substitution rate matrix for site  $i$ .

Estimating  $n$  additional parameters  $r_1, \dots, r_n$  is not feasible.

Instead estimate one meta-parameter  $\alpha$  and assume  $\Gamma$ -prior with shape parameter  $\alpha$  for all  $r_i$ .



The  $\Gamma$  distribution has another parameter, the scale parameter  $\beta$ . The expectation value of the  $\Gamma$  distribution is  $\alpha \cdot \beta$ .

We always assume  $\beta = 1/\alpha$ , such that

$$\mathbb{E}r_i = 1 \text{ and } \mathbb{E}Q = \mathbb{E}R_i$$

Density of the  $\Gamma$ -distribution:

$$g_{\alpha,\beta}(x) := \frac{x^{\alpha-1} \cdot e^{-x/\beta}}{\beta^\alpha \cdot \Gamma(\alpha)},$$

with  $\Gamma(a) = \int_0^\infty x^{a-1} \cdot e^{-x} dx$

We use

$$g_\alpha(x) := g_{\alpha,1/\alpha}$$

To contribution of data column  $D_i$  to the Likelihood of a tree  $T$  is then

$$L_{D_i}(T) = \Pr_T(D_i) = \int_0^\infty \Pr(D_i | r_i = x) \cdot g_\alpha(x) dx.$$

For each fixed  $r_i = x$  we can efficiently compute  $\Pr(D_i | r_i = x)$  with the Felsenstein pruning algorithm. But not for all  $x$  from 0 to  $\infty$ .

Idea: compute  $\Pr(D_i | r_i = x_j)$  for some  $x_j$  and approximate

$$\Pr(D_i) = \int_0^\infty \Pr(D_i | r_i = x) \cdot g_\alpha(x) dx \approx \sum_{j=1}^k w_j \cdot \Pr(D_i | r_i = x_j).$$

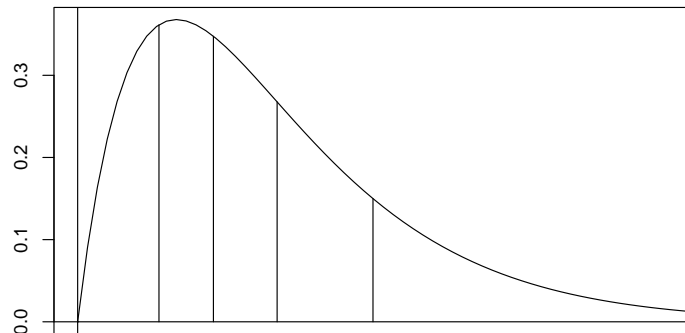
What are good choices for  $w_1, \dots, w_k$  and  $x_1, \dots, x_k$ ?

### Method of Yang (1994)

Divide  $[0, \infty]$  into  $k$  sections  $[a, b]$  of equal probability

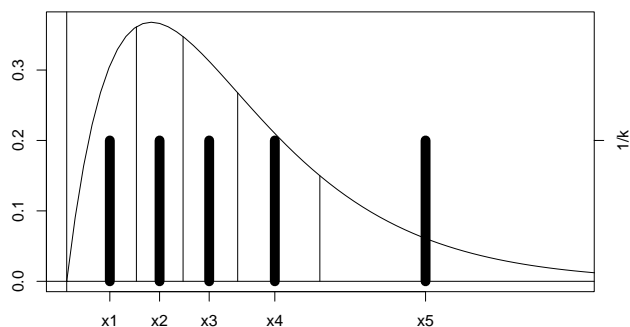
$$\int_a^b g_\alpha(x) dx = 1/k$$

, e.g. for  $k = 5$ :



Then,  $x_j$  is the expectation value of the  $\Gamma$  distribution conditioned of being in the  $j$ th section, i.e. the center of gravity of the area under the density.

All  $w_j$  are  $1/k$ .



Extension of the  $\Gamma$ -model: A proportion  $p$  of the sites is invariable.

### Alexis Stamatakis' CAT approximation

The "CAT model" provided by RAxML can be seen as an approximation to the discretized  $\Gamma$ -model.

- sites belong to a few different categories
- each category has its own rate acceleration factor that must be estimated
- ML estimate for each site to which category it belongs
- instead of marginalizing over all categories only use ML categories for likelihood computation
- Assignments of positions to categories are part of the parameter space and must be updated during ML optimization
- recommended if more than 50 taxa

Note: There is a completely different substitution model also called CAT in Lartillot and Philippe's program PhyloBayes.

## 5.7 Convergence into equilibrium

Let  $X = (X_1, X_2, \dots)$  or  $(X_t)_{t \in \mathbb{R}_{\leq 0}}$  be a Markov chain with finite state space  $\mathcal{Z}$  and transition probabilities  $P_{x \rightarrow y}(t)$  for  $t \in \mathbb{N}$  or  $t \in \mathbb{R}_{\geq 0}$ .

The transition dynamics  $P$  is **irreducible**, if

$$\forall x, y \in \mathcal{Z} \exists t : P_{x \rightarrow y}(t) > 0.$$

In the discrete-time case,  $P$  is **periodic**, if

$$\exists z \in \mathcal{Z}, k > 1 \forall n \in \mathbb{N} \setminus \{k, 2k, 3k, \dots\} P_{z \rightarrow z}(n) = 0$$

Otherwise,  $P$  is called **aperiodic**.

**Theorem 5** Each aperiodic irreducible transition dynamics  $P$  on a finite state space  $\mathcal{Z}$  has one and only one stationary distribution  $(\pi_z)_{z \in \mathcal{Z}}$ , i.e.

$$\forall z \in \mathcal{Z} \quad \pi_z = \sum_{x \in \mathcal{Z}} \pi_x \cdot P_{x \rightarrow z},$$

and converges against this distribution in the sense that

$$\forall x, z \quad \lim_{t \rightarrow \infty} P_{x \rightarrow z}(t) = \pi_z.$$

Sketch of proof of convergence: Start two Markov chains  $X$  and  $Y$  with transition matrix  $P$ , one with  $X_1$  in  $x$  and one with  $Y_1$  taken from the stationary distribution. When they meet in some step  $k$ , i.e. if  $X_k = Y_k$ , couple them:  $X_j = Y_j$  for all  $j > k$ . If  $P$  is irreducible and aperiodic, and the probability  $q_k$  that  $X$  and  $Y$  do not meet before step  $k$  converges to 0, and

$$|\Pr(X_j = z) - \pi_z| = |\Pr(Y_j = z) - \Pr(X_j = z)| = |\Pr(Y_j = z, X_j = Y_j) + \Pr(Y_j = z, X_j \neq Y_j) - \Pr(X_j = z, X_j = Y_j) - \Pr(X_j = z, X_j \neq Y_j)| = |\Pr(Y_j = z, X_j \neq Y_j) - \Pr(X_j = z, X_j \neq Y_j)| \leq \max\{\Pr(Y_j = z, X_j \neq Y_j), \Pr(X_j = z, X_j \neq Y_j)\} \leq q_j \rightarrow 0.$$

A Markov chain with transition matrix  $P$  and stationary distribution  $(\pi_z)_{z \in \mathcal{Z}}$  is **reversible** if

$$\forall_{z,y \in \mathcal{Z}} : \pi_z \cdot P_{z \rightarrow y}(t) = \pi_y \cdot P_{y \rightarrow z}(t).$$

(“detailed-balance condition”)

Note: the detailed-balance condition already implies that  $(\pi_z)_{z \in \mathcal{Z}}$  is a stationary distribution of  $P$ .

The evolutionary dynamics described by Jukes-Cantor, F81, F84, HKY, GTR or PAM matrices are reversible. If we assume reversibility and no molecular clock, the likelihood does not depend on the position of the root in the tree topology.

If the root divides a branch of length  $s + t$  into sections of length  $s$  and  $t$ , reversibility implies that the probability stays the same if we move the root into one of the nodes:

$$\begin{aligned} \sum_z \pi_z \cdot P_{z \rightarrow x}(s) \cdot P_{z \rightarrow y}(t) &= \sum_z \pi_x \cdot P_{x \rightarrow z}(s) \cdot P_{z \rightarrow y}(t) \\ &= \pi_x \cdot P_{x \rightarrow y}(s+t) \\ &= \pi_y \cdot P_{y \rightarrow x}(s+t) \end{aligned}$$

## 6 Bayesian phylogeny reconstruction and MCMC

### 6.1 Principles of Bayesian statistics

In Bayesian statistics, also parameters are equipped with probabilities.

For phylogenetic trees  $T$ :

**prior probability distribution:**  $P(T)$  is the probability density of the tree  $T$  disregarding the data, e.g. we could a priori assume a uniform probability density for all trees up to a certain total branch length.

**posterior probability distribution**  $P(T|D)$  is the conditional probability density of the tree  $T$ , given the data  $D$ .

Bayes-Formula:

$$P(T|D) = \frac{P(T, D)}{\Pr(D)} = \frac{\Pr(D|T) \cdot P(T)}{\int_{T'} \Pr(D|T') \cdot P(T') dT'}$$

Computing

$$P(T|D) = \frac{\Pr(D|T) \cdot P(T)}{\int_{T'} \Pr(D|T') \cdot P(T') dT'}$$

is not trivial. We can compute  $\Pr(D|T) = \Pr_T(D) = L_D(T)$  by Felsenstein pruning and  $P(T)$  is defined by our prior distribution, but integrating over all trees is difficult.

What we can compute is the ratio of the probabilities of two candidate trees  $T_A$  and  $T_B$ :

$$\frac{P(T_A|D)}{P(T_B|D)} = \frac{\frac{\Pr(D|T_A) \cdot P(T_A)}{\int_{T'} \Pr(D|T') \cdot P(T') dT'}}{\frac{\Pr(D|T_B) \cdot P(T_B)}{\int_{T'} \Pr(D|T') \cdot P(T') dT'}} = \frac{\Pr(D|T_A) \cdot P(T_A)}{\Pr(D|T_B) \cdot P(T_B)}$$

## 6.2 MCMC sampling

We are not just interested in finding the *maximum a-posteriori* (MAP) tree

$$\arg \max_T P(T|D),$$

but, very much in the spirit of Bayesian statistics, to sample trees from the posterior distribution, that is, to generate a set of (approximately) independent random trees  $T_1, T_2, \dots, T_n$  according to the probability distribution given by  $P(T|D)$ . This will allow us not only to infer the phylogeny but also to assess the uncertainty of this inference.

Idea: Simulate a Markov chain on the space of trees with stationary distribution  $P(T|D)$  and let it converge.

How can we do that if we can only compute ratios  $\frac{P(T_A|D)}{P(T_B|D)}$  for given trees  $T_A$  and  $T_B$ ?

Given the probability distribution  $\Pr(\cdot|D)$ , how can we construct a Markov chain that converges against it?

One possibility: **Metropolis-Hastings**

Given current state  $X_i = x$  propose  $y$  with Prob.  $Q(x \rightarrow y)$

Accept proposal  $X_{i+1} := y$  with probability

$$\min \left\{ 1, \frac{Q(y \rightarrow x) \cdot \Pr(y | D)}{Q(x \rightarrow y) \cdot \Pr(x | D)} \right\}$$

otherwise  $X_{i+1} := X_i$

(All this also works with continuous state space, with some probabilities replaced by densities.)

### Why Metropolis-Hastings works

Let's assume that  $\frac{Q(y \rightarrow x) \cdot \Pr(y | D)}{Q(x \rightarrow y) \cdot \Pr(x | D)} \leq 1$ . (Otherwise swap  $x$  and  $y$  in the following argument). Then, if we start in  $x$ , the probability  $\Pr(x \rightarrow y)$  to move to  $y$  (i.e. first propose and then accept this) is

$$Q(x \rightarrow y) \cdot \frac{Q(y \rightarrow x) \cdot \Pr(y | D)}{Q(x \rightarrow y) \cdot \Pr(x | D)} = Q(y \rightarrow x) \frac{\Pr(y | D)}{\Pr(x | D)}$$

If we start in  $y$ , the probability  $\Pr(y \rightarrow x)$  to move to  $x$  is

$$Q(y \rightarrow x) \cdot 1,$$

since our assumption implies  $\frac{Q(x \rightarrow y) \cdot \Pr(x | D)}{Q(y \rightarrow x) \cdot \Pr(y | D)} \geq 1$ .

This implies that the reversibility condition

$$\Pr(x | D) \cdot \Pr(x \rightarrow y) = \Pr(y | D) \cdot \Pr(y \rightarrow x)$$

is fulfilled. This implies that  $\Pr(\cdot | D)$  is an equilibrium of the Markov chain that we have just constructed, and the latter will converge against it. (let's watch a simulation in R)

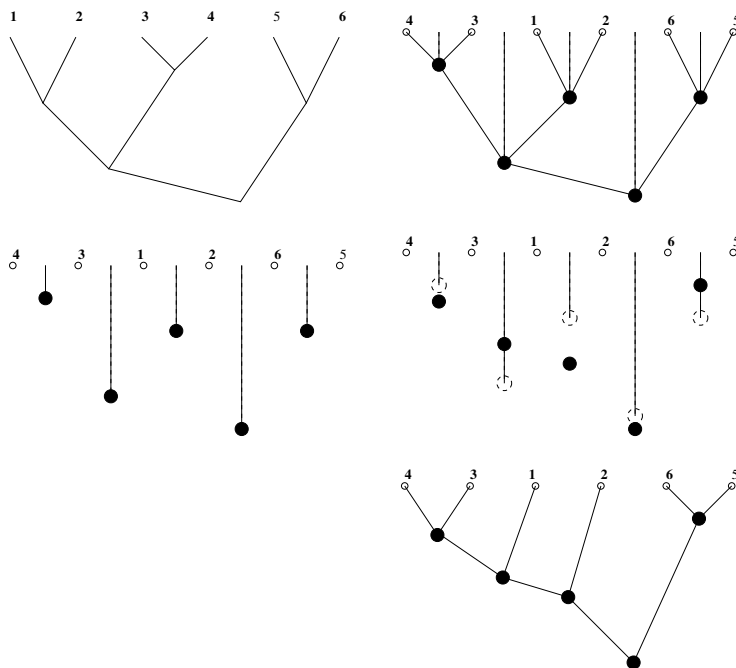
### Applying Metropolis-Hastings

- You are never in equilibrium (your target distribution), but you can get close if you run enough steps.
- You can take more than one sample from the same chain, but you should run enough steps between the sampling steps to make the sampled objects only weakly dependent.
- Your initial state may be “far from equilibrium” (i.e. very improbable). So you should run the chain long enough before you start sampling (“burn-in”).
- Launch many independent MCMC runs with different starting points and check whether they lead to the same results.

**Mau, Newton, Larget 1996**

Seminal paper on MCMC for phylogenies; propose a proposal chain for ultrametric trees.

1. Draw the tree in the plane.
2. In each internal node rotate subtrees with probability 1/2.
3. Remove edges from drawing.
4. Shift each internal node in time by a random amount.
5. Reconstruct edges from modified time points of nodes.



Other software packages use more common tree modifications like NNI, SPR and TBI.

Examples of software for Bayesian sampling:

**MrBayes** <http://mrbayes.csit.fsu.edu/>

**BEAST** [http://beast.bio.ed.ac.uk/Main\\_Page](http://beast.bio.ed.ac.uk/Main_Page)

**PhyloBayes** <http://www.atgc-montpellier.fr/phylobayes/binaries.php>

**BAl-i-Phy** <http://www.biomath.ucla.edu/msuchard/bali-phy/>

**TreeTime** [http://evol.bio.lmu.de/\\_statgen/software/treetime/](http://evol.bio.lmu.de/_statgen/software/treetime/)

**(MC)<sup>3</sup>=MCMCMC**

=Metropolis-Coupled MCMC= MCMC with “heated chains”.

If  $\beta_i \in (0, 1]$  is heat parameter for chain  $i$ , then chain  $i$  samples from distribution  $p^{\beta_i} : x \mapsto p^{\beta_i}(x) \cdot \text{const}$ , with  $\beta_1 = 1$ .

The usual MH acceptance prob. for chain  $i$  is

$$\min \left\{ 1, \frac{p(y)^{\beta_i}}{p(x)^{\beta_i}} \cdot \frac{Q_{y \rightarrow x}}{Q_{x \rightarrow y}} \right\}.$$

Sometimes a swap between the current state  $x_i$  of chain  $i$  and the current state  $x_j$  of chain  $j$  is proposed. The acceptance with probability

$$\min \left\{ 1, \frac{p(x_i)^{\beta_i}}{p(x_j)^{\beta_i}} \cdot \frac{p(x_j)^{\beta_j}}{p(x_i)^{\beta_j}} \right\}$$

fulfills the requirements of both chains (check this!).

Most programs for Bayesian phylogeny inference can also estimate parameters of the substitution model. Combine the estimation of trees with the estimation of divergence times or even alignments.

*Gibbs sampling* is applied to combine Bayesian estimations for different kinds of parameters.

### Gibbs sampling

Assume we want to sample from a joint distribution  $\Pr(A = a, B = b)$  of two random variables, and for each pair of possible values  $(a, b)$  for  $(A, B)$  we have Markov chains with transition probabilities  $P_{b \rightarrow b'}^{(A=a)}$  and  $P_{a \rightarrow a'}^{(B=b)}$  that converge against  $\Pr(B = b|A = a)$  and  $\Pr(A = a|B = b)$ .

Then, any Markov chain with transition law

$$P_{(a,b) \rightarrow (a',b')} = \begin{cases} \frac{1}{2}P_{a \rightarrow a'}^{(B=b)} + \frac{1}{2}P_{b \rightarrow b'}^{(A=a)} & \text{if } a = a' \text{ and } b = b' \\ \frac{1}{2}P_{a \rightarrow a'}^{(B=b)} & \text{if } a \neq a' \text{ and } b = b' \\ \frac{1}{2}P_{b \rightarrow b'}^{(A=a)} & \text{if } a = a' \text{ and } b \neq b' \\ 0 & \text{else} \end{cases}$$

## 6.3 Checking convergence of MCMC

### Effective Sampling Size (ESS)

Assume that we want to estimate the expectation value  $\mu$  of a distribution by taking the mean  $\bar{X}$  of  $n$  independent draws  $X_1, X_2, \dots, X_n$  from the distribution. Then,

$$\begin{aligned} \mathbb{E}\bar{X} &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}X_i = \mu \\ \text{var}(\bar{X}) &= \text{var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \text{var}(X_1). \end{aligned}$$

If we instead use  $m$  *correlated* draws  $Y_1, Y_2, \dots, Y_m$  from the same distribution, then

$$\begin{aligned} \mathbb{E}\bar{Y} &= \frac{1}{m} \sum_{i=1}^m \mathbb{E}Y_i = \mu \\ \text{var}(\bar{Y}) &= \text{var}\left(\frac{1}{m} \sum_{i=1}^m Y_i\right) = \frac{1}{m} \text{var}(Y_1) + \frac{2}{m^2} \sum_{i=1}^m \sum_{j=i+1}^m \text{cov}(Y_i, Y_j). \end{aligned}$$

### Effective Sampling Size (ESS)

With the autocorrelation  $\rho_k = \text{cor}(Y_i, Y_{i-k}) = \text{cor}(Y_i, Y_{i-k})/\text{var}(Y_i)$ ,  $\bar{Y}$  has (approximately) the same variance as  $\bar{X}$ , if

$$n = \frac{m}{1 + 2 \cdot \sum_{k=1}^{\infty} \rho_k}.$$

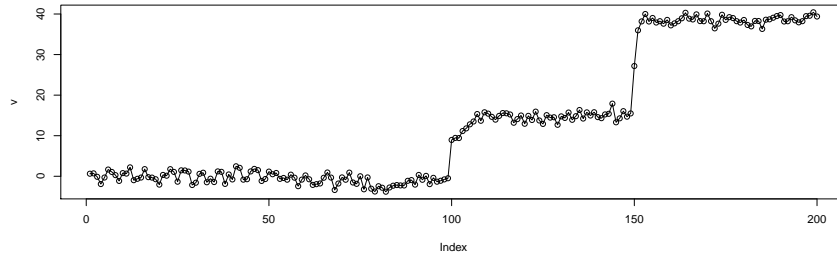
Therefore, we estimate the Effective Sample Size by

$$ESS = \frac{m}{1 + 2 \cdot \sum_{k=1}^{\infty} \hat{\rho}_k},$$

where  $\hat{\rho}_k$  is an estimation of the autocorrelation  $\rho_k := \text{cor}(Y_i, Y_{i-k})$ .



Problem: ESS may be too optimistic because correlation may be underestimated.



estimated effective sample sizes:

range	:	ess
1-90	:	7.88
110-140	:	31.00
160-200	:	28.77
1-200	:	1.53

Ways to check convergence of MCMC

- ESS
- visually inspect paths of log likelihood and parameter estimates
- start many MCMC runs with different start values and check whether they appear to converge against the same distribution

## 6.4 Interpretation of posterior probabilities and robustness

If the prior is correctly chosen and the model assumptions are fulfilled, the posterior probability of a tree topology should be the probability that the topology is correct. This is confirmed for trees with six taxa in a simulation study in:

## References

[HR04] J.P. Huelsenbeck, B. Rannala (2004) Frequentist Properties of Bayesian Posterior Probabilities of Phylogenetic Trees Under Simple and Complex Substitution Models. *Syst. Biol.* **53**(6):904–913.

However, when a the model chosen for the substitution process is too simple (e.g. neglecting rate heterogeneity), the estimated posterior probabilities can be over-optimistic. Using a model that is more complex than necessary, may lead to just slightly conservative estimates of posterior probabilities. Recommendation: If you are not sure, rather use the more complex substitution model.

## References

[YR05] Z. Yang, B. Rannala (2005) Branch-Length Prior Influences Bayesian Posterior Probability of Phylogeny *Syst. Biol.* *54*(3):455–470

simulate rooted ultrametric trees with three tips and different priors for lengths of inner and outer branches. Compute posterior probabilities for the three possible topologies with various priors for tree lengths.

- MAP estimates are robust against misspecification of prior.
- High posteriors are underestimated and low posteriors are overestimated if prior favors very short internal edges.
- High posteriors are overestimated and low posteriors are underestimated if priors for internal edge lengths are flat.

Note: flat priors are sometimes called “uninformative”, but this is misleading, and in Yang and Rannala’s study these priors were most problematic!

To decrease the risk of too optimistic posteriors for tree topologies when the substitution process is inappropriate,

## References

[Y08] Z. Yang (2008) Empirical evaluation of a prior for Bayesian phylogenetic inference *Phil. Trans. R. Soc. B* **363**: 4031–4039

recommends using priors favoring shorter branch lengths if the input alignment is long.

### Star-tree paradox

If the inner branch of a rooted 3-taxa tree is extremely short, or even non-existing, and the Bayesian method takes only binary trees into account with “liberal” priors for the branch lengths, it will often assign a high posterior probability to one of the three tree topologies, and with probability  $\approx 2/3$  it will be a wrong one.

This is related to the *fair-coin paradox* and *Lindley’s paradox*, which we will discuss in the context of Bayesian model selection.

## References

[MV05] E. Mossel, E. Vigoda (2005) Phylogenetic MCMC Algorithms Are Misleading on Mixtures of Trees. *Science* **309**: 2207–2209

point out that Bayesian MCMC procedures may assign a high posterior probability to a tree when the data is actually a mixture of data from two different trees. See also

## References

[RLH+06] F. Ronquist, B. Larget, J.P. Huelsenbeck, J.B. Kadane, D. Simon, P. van der Mark (2006) Comment on “Phylogenetic MCMC Algorithms Are Misleading on Mixtures of Trees” *Science* **312**:367a

[MV06] E. Mossel, E. Vigoda (2006) Response to Comment on “Phylogenetic MCMC Algorithms Are Misleading on Mixtures of Trees” *Science* **312**:367b

## 7 Bootstrapping

### 7.1 The concept of bootstrapping

Assume a panmictic Hardy-Weinberg population and a locus in equilibrium with genotypes  $MM$ ,  $MN$ , and  $NN$ . This means, the frequencies of these genotypes are  $(1 - \theta)^2$ ,  $2\theta(1 - \theta)$ , and  $\theta^2$ , where  $\theta$  is the frequency of allele  $N$ .

Assume the following observations:

$MM$	$MN$	$NN$	total
342	500	187	1029
$X$	$Y$	$Z$	

(Example taken from Rice (1995) *Mathematical Statistics and Data Analysis*. Duxbury press.)

We estimate  $\theta$  by  $\hat{\theta} = \frac{2Z+Y}{2(X+Y+Z)} = 0.4247$ . How accurate is this estimation?

Simulate 1000 datasets, each consisting of 1029 individuals drawn from a Hardy-Weinberg population with frequency 0.4247 of allele  $N$ .

Let  $\theta_1^*, \theta_2^*, \dots, \theta_{1000}^*$  be the estimates of  $\theta$  from the 1000 datasets. We can then estimate the standard deviation of our estimator  $\hat{\theta}$  by

$$\sigma_{\hat{\theta}} \approx \sqrt{\frac{\sum_i (\theta_i^* - \hat{\theta})^2}{999}}$$

Bootstrapping is a general approach in statistics that is often used to assess the accuracy of an estimator.

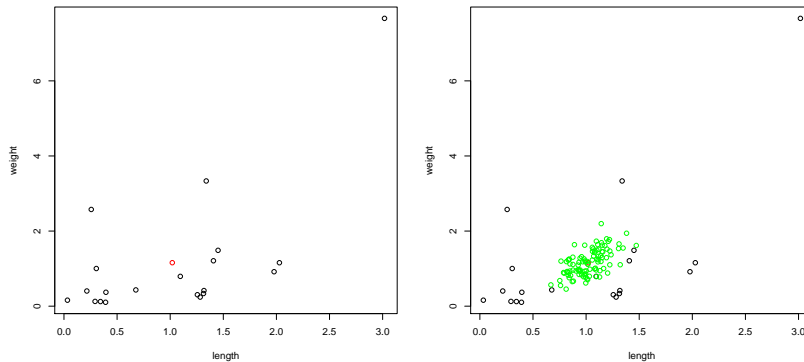
It is based on the following idea: If we estimate a parameter  $\theta$  by  $\hat{\theta}$ , we can check the accuracy of the estimation method with simulated data.

Problem: We do not know the true value of  $\theta$  but need a value for the simulations.

idea: We pull ourselves up by our own bootstraps by using  $\hat{\theta}$  for the simulations and assume that the difference  $\hat{\theta} - \theta^*$ , where  $\theta^*$  is the estimation from the simulated data, has a similar distribution as  $\theta - \hat{\theta}$ :

$$\mathcal{L}(\theta - \hat{\theta}) \approx \mathcal{L}(\hat{\theta} - \theta^*)$$

Since we use the parameter and assumptions about its distribution, this is called *parametric bootstrap*. In the next example we use *non-parametric bootstrapping*, which means that we just use the original data to simulate new data.

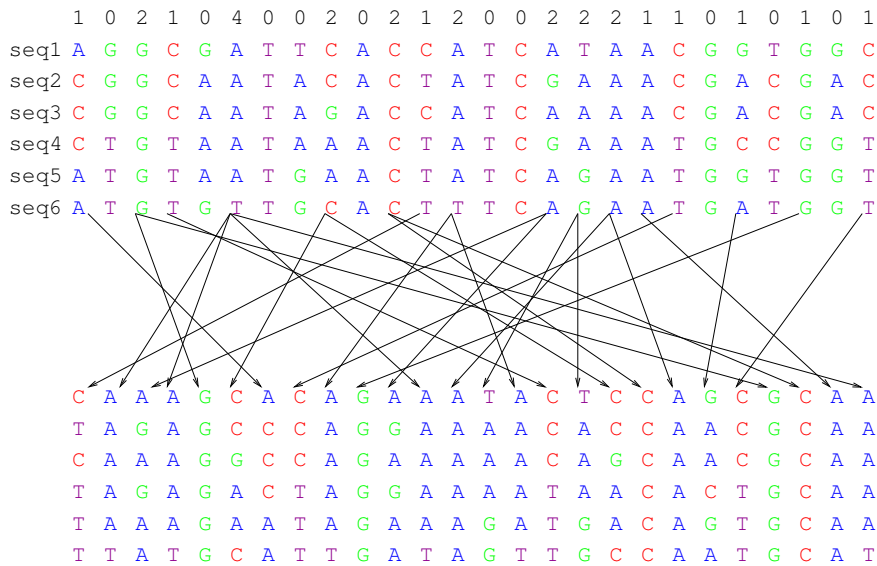


We have caught 20 fishes from a lake and want to estimate the distribution of size and weight in the population by the sample means. How accurate is this estimation? Idea: simulate sampling from a population by putting the 20 fishes into a pond and take a sample of size 20. To avoid getting precisely the same sample, *sample with replacement*. Compute the mean length and weight from the “bootstrap sample”. Repeat this procedure 1000 times. The 1000 pairs of means can be used for bias correction and to estimate the variance of the estimator.

$$\text{Bias correction: } \hat{\theta} - (\overline{\theta^*} - \hat{\theta}) = 2\hat{\theta} - \overline{\theta^*}$$

## 7.2 Bootstrap for phylogenetic trees

### non-parametric bootstrap of an alignment



A bootstrap alignment has the same length as the original alignment. It consists of columns that were randomly drawn from the original alignment with replacement. To the bootstrap alignment we apply the same phylogeny reconstruction method as for the original alignment.

We repeat this many times and thus get many bootstrap trees. We label each branch of our originally reconstructed tree by the percentages of bootstrap trees that have this branch. These bootstrap values are supposed to give an impression of how reliable the branches are.

Alternatives to non-parametric bootstrap:

**Jackknife:** Create shorter alignments, e.g. 90%, by sampling without replacement. Like in non-parametric bootstrapping, the bootstrap dataset is slightly less informative than the original data.

**Parametric Bootstrap:** Use the estimated tree and substitution rates estimated along with the tree to simulate new data. (Disadvantage: does not take uncertainty about the substitution model into account.)

### 7.3 How can we interpret the bootstrap values?

There are at least three different interpretations of the bootstrap values of tree branches:

1. posterior probability of the branch
2. measures of repeatability
3. confidence levels for the existence of the branch

None of these interpretations is perfect.

Are bootstrap values posterior probabilities?

Rather not, because posterior probabilities depend on the prior, and the bootstrap values do not (at least if a non-Bayesian method was used for tree reconstruction).

Do bootstrap values measure repeatability?

This is the original interpretation of Felsenstein, who first proposed bootstrapping for phylogenetic trees. However, the bootstrap value can only be an approximative measure because the bootstrap sample is slightly less informative than the original sample. The question is also what repeatability would actually mean? If the analysis is repeated with different data, variations between loci may play a role, which is not incorporated in bootstrapping.

Are bootstrap values confidence levels?

If a branch has a bootstrap value 97% and this is interpreted as confidence level, then this means the following: Under the null hypothesis that the branch is actually not there or has length 0, the probability

of getting a bootstrap support of 97% is  $100\% - 97\% = 3\%$ . This means: Among all branches that appear in the estimated trees but are actually wrong, only 3% get such a high bootstrap level.

It has been conjectured that bootstrap values underestimate confidence because bootstrap datasets are less informative than the original dataset. However, this argument disregards that the bootstrap result  $\theta^*$  does not need to be an approximation for  $\hat{\theta}$ , but  $\theta^* - \hat{\theta}$  should be an approximation for  $\hat{\theta} - \theta$ .

## References

[EHH96] B. Efron, E. Halloran, S. Holmes (1996) Bootstrap confidence levels for phylogenetic trees. *Proc. Nat. Acad. Sci. U.S.A.* **93(13)**:429–434

show that bootstrap values can either over- or underestimate confidence, but are at least first-order approximations of confidence values. They propose a meta-bootstrap procedure to correct the over- or underestimation for each branch.

## 8 Tests for trees and branches

### 8.1 The Kishino-Hasegawa (KH) test

The KH test compares two given trees. The null hypothesis is that differences in their likelihoods are only due to “sampling error”, i.e. the mutations that randomly occurred at the sites in our dataset. Several versions of the KH test exist, one of them is as follows:

- Given an alignment of length  $S$  let for each  $k \leq S$  be  $\ell_1^{(k)}$  and  $\ell_2^{(k)}$  the log-likelihoods of the two trees for the  $k$ -th column of the alignment.
- define  $\delta_k = \ell_1^{(k)} - \ell_2^{(k)}$
- estimate the variance of all  $\delta^k$  by  $\hat{\sigma}^2 = \frac{\sum_k (\delta_k - \bar{\delta})^2}{S-1}$ , where  $\bar{\delta}$  is the mean over all  $\delta_k$ .
- Under the null hypothesis (and model assumptions like independence of sites etc.), the log likelihood-ratio  $\ell_1 - \ell_2$  is normally distributed with mean 0 and variance  $S \cdot \sigma^2$ .

Hence, reject the null hypothesis on the 5% level if  $|\ell_1 - \ell_2| > 1.96 \cdot \sqrt{S\hat{\sigma}}$

(other variants of the test use log likelihood-ratios of bootstrapped trees instead of site-wise log likelihood-ratios)

Note that the selection of trees to be tested must be independent from the data that is used in the KH test!

If one of the trees has been selected because of its high likelihood for this dataset, the other tree will be rejected too often!

To apply the KH test to more than two trees, some multiple-testing correction is needed.

### 8.2 The Shimodaira-Hasegawa (SH) test

Assume that a set of tree topologies is given that includes the true topology. Again, the choice of the set of topologies must be independent of the data. The null hypothesis is that differences in the likelihoods of the trees are only due to “sampling error”.

1. Make  $R$  bootstrap samples from the  $S$  sites and compute the log likelihood  $\ell_{t,r}$  for each tree  $t$  in the set and each bootstrapped data set  $r$ .
2.  $\tilde{R}_{t,r} := \ell_{t,r} - \frac{1}{R} \sum_{k=1}^R \ell_{t,k}$
3.  $D_{t,r} := \max_p \tilde{R}_{p,r} - \tilde{R}_{t,r}$

4. A tree  $t$  can be rejected on the 5% level if at most 5% of its  $D_{t,r}$  values are higher than the difference between the maximum log-likelihood and the log-likelihood  $\ell_t$  of  $t$ .

Note that this has a built-in multiple-testing correction!

### 8.3 The SOWH test

## References

- [GAR00] N. Goldman, J.P. Anderson, A.G. Rodrigo (2000) Likelihood-Based Tests of Topologies in Phylogenetics *Syst. Biol.* **49(4)**: 652–670
- [SOWH] D.L. Swofford, G.J. Olsen, P.J. Waddell, D.M. Hillis (1996) Phylogenetic inference in: D.M. Hillis, C. Moritz, B.K. Mabe (eds.) *Molecular Systematics*, Sinauer.

To test whether a tree  $T_1$  can be rejected ( $H_0$ : “ $T_1$  is the true tree”), use as a test statistic the difference  $\delta = \ell_{ML} - \ell_1$  between the maximum log likelihood  $\ell_{ML}$  and the log likelihood  $\ell_1$  of  $T_1$ .

Simulate many datasets  $d$  by parameteric bootstrapping using  $T_1$  and the corresponding estimates of all parameters (mutation rates, branch lengths etc.).

Let  $\ell_{1,d}$  be the likelihood of  $T_1$  based on bootstrap data set  $d$  with new estimations for all parameters, and let  $\ell_{ML,d}$  be the same maximized over all tree topologies.

Use all  $\delta_d = \ell_{ML,d} - \ell_{1,d}$  (for all  $d$ ) to estimate the distribution of the test statistic  $\delta$  under the null hypothesis that  $T_1$  is correct.

Reject  $T_1$  on the 5% level if less than 5% of the  $\delta_d$  are larger than  $\delta$ .

## References

- [B02] T.R. Buckley (2002) Model Misspecification and Probabilistic Tests of Topology: Evidence from Empirical Data Sets *Syst. Biol.* **51(3)**: 509–523

Shows examples where SOWH test and posterior probabilities falsely reject too many trees because of using the wrong substitution models. The SH test does not have this problem and rather tends to be too conservative.

Uses real data with phylogeny more or less well known.

**Advantage:** Realistic because all substitution model used in simulation study are somehow idealized.

**Drawbacks:** Only a few such datasets are available and results may not be representative. In principle, the assumed phylogenies could still be erroneous.

### 8.4 An approximate Likelihood-Ratio Test (aLRT)

## References

- [AG06] M. Anisimova, O. Gascuel (2006) Approximate Likelihood-Ratio Test for Branches: A Fast, Accurate, and Powerful Alternative *Syst. Biol.* **55(4)**: 539–552

We want to show significance of a particular branch in the tree, i.e. the null hypothesis is that this tree has length 0. However, we assume that with any other respect, the topology of the tree is true.

A likelihood-ratio test is applied to check whether the model with the tree being longer than 0 fit the data significantly better than a model where the branch length is restricted to 0.

Thus, we have to approximate the distribution of the (log) likelihood-ratio under the null hypothesis.

Ususally, if we have a model  $M_1$  with  $n - d$  parameters nested in a model  $M_2$  with  $n$  parameters, then under the null-hypothesis that the data come from the more simple model  $M_1$ , the double log likelihood ratio is approximately chisquare-distributed with  $d$  degrees of freedom,

$$\mathcal{L}_{M_1} \left( 2 \cdot \log \frac{L_D(M_2)}{L_D(M_1)} \right) = \mathcal{L}_{M_1} (2 \cdot (\log L_D(M_2) - \log L_D(M_1))) \approx \chi_d^2,$$

where the likelihood of a model  $L_D(M_i)$  is maximum likelihood obtained by optimization over all parameters of the model.

This only works, if the model  $M_1$  is in the inner of the model  $M_2$ . In our case, the null hypothesis is at the boundary of the more general model, because the branch length 0 is on the boundary of the set of allowed branch lengths. Therefore, the following correction proposed: The distribution of  $2 \cdot (\log L_D(M_2) - \log L_D(M_1))$  is approximated by a distribution that puts weight 0.5 on 0 and half of the density of  $\chi_1^2$  on all positive values.

Anisimova, Gascuel (2006): Let  $\ell_1$  be the log likelihood of the ML tree,  $\ell_0$  that of the topology with the length of the focal branch removed, and  $\ell_2 > \ell_3$  the log likelihoods of the two topologies where the focal branch is removed in an NNI step and (see Figure 1 in Anisimova, Gascuel (2006))[1.5ex] For more robustness,  $2(\ell_1 - \ell_2) \leq 2(\ell_1 - \ell_0)$  is used as a test statistic. (Maybe the idea is that the null hypothesis should be that one of the other fully resolved trees is right.)[1.5ex] The likelihood of a topology is the maximum likelihood of a tree with this topology. Thus, each value  $\ell_0, \ell_2, \ell_3$  needs own optimization of all branch lengths. Here, Anisimova and Gascuel use an approximation by optimizing only the four neighboring branches of the focal branch and its alternative branch in the case of  $\ell_2$  and  $\ell_3$ . [1.5ex] If the null hypothesis is true, any of the three possible fully resolved topologies can get the highest likelihood. Therefore, a multipl e-testing correction is needed. The Bonferroni correction is applied, which means that the  $\alpha$ -level is replaced by its third.

Anisimova and Gascuel conclude from simulations that

- Approximate likelihood-ratio test (aLRT, i.e. with optimization over only five branches) has accuracy and power similar to standard LRT.
- aLRT is robust against mild model misspecifications.
- aLRT was slightly more accurate w.r.t. 5% type I error than ML bootstrap.
- In contrast to common belief, bootstrap was a bit too liberal, i.e. its type I error rate was higher than the significance level.
- Bayesian methods were a bit too conservative in this simulation study.

## 9 Model selection

### 9.1 Concepts: AIC, hLRT, BIC, DT, Model averaging, and bootstrap again

#### AIC

The likelihood of a model  $M$ ,

$$L_D(M) = \max_{\theta \in M} L_D(\theta) = \max_{\theta} \Pr_{M,\theta}(D)$$

tells us how well  $M$  fits the data  $D$ . The more parameter dimensions  $d$  (i.e.  $\theta = (\theta_1, \theta_2, \dots, \theta_d)$ ) the higher the likelihood and the higher the risk of *overfitting*!

Under certain assumptions (with normal distributions, not phylogenies), the error of future predictions in terms of Kullback-Leibler-Information can be estimated by Akaike's Information Criterion:

$$\text{AIC} = -2 \cdot \log L_D(M) + 2 \cdot d.$$

One approach: use the model of lowest AIC.

### Model selection via LRT

If we have a model  $M_1$  with  $n - d$  parameters nested in a model  $M_2$  with  $n$  parameters, then under the null-hypothesis that the data come from the more simple model  $M_1$ , the double log likelihood ratio is under certain conditions approximately chisquare-distributed with  $d$  degrees of freedom,

$$\mathcal{L}_{M_1} \left( 2 \cdot \log \frac{L_D(M_2)}{L_D(M_1)} \right) = \mathcal{L}_{M_1} (2 \cdot (\log L_D(M_2) - \log L_D(M_1))) \approx \chi_d^2,$$

where the likelihood of a model  $L_D(M_i)$  is maximum likelihood obtained by optimization over all parameters of the model.

In cases where the  $\chi_d^2$  approximation is dubious (e.g. when the models are not nested) one can simulate the likelihood ratio distribution under the null hypothesis.

One approach of model selection is to accept the more complex model only if the simpler model is significantly violated.

### Problems of this LRT approach

- Model selection is different from the original idea of testing. If a test does not show significance, one cannot conclude anything, and especially not that the null hypothesis (the simpler model) is favorable.
- One can in principle apply this to a hierarchy of nested models, but the result will depend on which intermediate steps are allowed.

### Bayesian model selection

Each model  $M_i$  has a prior probability  $\Pr(M_i)$ . Its posterior probability is then

$$\Pr(M_i|D) = \frac{\Pr(D|M_i) \cdot \Pr(M_i)}{\sum_j \Pr(D|M_j) \cdot \Pr(M_j)}$$

with

$$\Pr(D|M_i) = \int_{\theta} \Pr(D|M_i, \theta) \cdot \Pr(\theta|M_i) d\theta.$$

Note the difference between  $\Pr(D|M_i)$ , where we integrate over  $\theta$ , and  $L_D(M_i)$  where we maximize over  $\theta$ ! The sum over all models in the denominator above cancels if we compare two models by taking the ratios of their posteriors:

$$\frac{\Pr(M_1|D)}{\Pr(M_2|D)} = \frac{\Pr(D|M_1)}{\Pr(D|M_2)} \cdot \frac{\Pr(M_1)}{\Pr(M_2)}$$

The fraction  $\Pr(D|M_1)/\Pr(D|M_2)$  is called the *Bayes factor* of the models  $M_1$  and  $M_2$ .

To avoid the priors of the models we use the Bayes factors rather than the posterior distributions to decide between models. If the Bayes factor  $\Pr(D|M_1)/\Pr(D|M_2)$  is larger than 1 we may favor  $M_1$  over  $M_2$ . The rule of thumb says that a Bayes factor between 1 and 3 is not worth mentioning, between 3 and 20 it indicates some evidence, between 20 and 150 strong evidence, and over 150 very strong evidence.

It is important to note that even if the priors  $\Pr(M_i)$  of the models do not matter, the priors  $\Pr_{M_i}(\theta)$  within the model may have a strong influence. An important difference between Bayesian parameter estimation and Bayesian model selection is that priors become less important for parameter estimation as more data is added. This is not the case in model selection, where priors for the model parameters will always have an important impact!

Some people find the following properties of posterior probabilities counter-intuitive:

**Lindley's paradox** In the limit of uninformative priors, the simplest model is always preferred.

**Star-tree paradox** If all internal nodes have length (almost) 0, there will often be a fully resolved tree with high posterior probability (deciding between topology can be considered as model selection).

**Fair-coin paradox** If a (almost) fair coin is tossed many times, but the models compared allow only for one or the other side to have probability larger than 0.5, it will often be the case that one of the two models has a high posterior probability.



### Computation of Bayes factors from MCMC runs

If  $\theta^{(1)}, \dots, \theta^{(m)}$  are (approximately) independent samples  $\Pr(\theta|D, M)$  we can compute  $\Pr(D|M)$  by importance sampling approximation:

$$\begin{aligned} \Pr(D|M) &= \frac{\Pr(D|M)}{\int_{\theta} \Pr(\theta|M) d\theta} = \frac{1}{\int_{\theta} \frac{\Pr(\theta|M)}{\Pr(D|M)} d\theta} \\ &\approx \frac{1}{\frac{1}{m} \sum_{\theta^{(i)}} \frac{\Pr(\theta^{(i)}|M)}{\Pr(D|M) \cdot \Pr(\theta^{(i)}|D, M)}} \\ &= \frac{m}{\sum_{\theta^{(i)}} \frac{\Pr(\theta^{(i)}|M)}{\Pr(D, \theta^{(i)}|M)}} = \frac{m}{\sum \frac{1}{\Pr(D|M, \theta^{(i)})}} \end{aligned}$$

(note that this harmonic mean estimator may be numerically instable.)

### BIC

For a model  $M$  with a  $d$ -dimensional parameter  $\theta$  and data  $D$  consisting of  $N$  independent samples, we can under certain conditions approximate

$$\log \Pr(D|M) \approx \log \Pr(D|M, \hat{\theta}) - \frac{d}{2} \cdot \log N$$

We call  $BIC(M) = -2 \cdot \log \Pr(D|M, \hat{\theta}) + d \log N$  the *Bayesian Information Criterion* or *Schwartz Criterion*, and favor models of low  $BIC$ . Moreover,

$$\frac{\Pr(D|M_1)}{\Pr(D|M_2)} \approx e^{(BIC(M_2) - BIC(M_1))/2}.$$

Minin, Abdo, Joyce, Sullivan (2003): “[...] rather than worry about the somewhat artificial criterion whether or not a model is correct, we will focus on the accuracy of the branch lengths estimated under various models”

- Assume that the unrooted phylogeny with  $k$  tips is known or use some initial guess.
- Candidate models:  $M_1, \dots, M_m$  with uniform prior  $\Pr(M_i) = 1/m$ .
- Branch lengths estimated with model  $M_i$ :

$$\begin{aligned} B_i &= (\widehat{B}_{i,1}, \dots, \widehat{B}_{i,2k-3}) \\ \|B_i - B_j\| &:= \sqrt{\sum_{\ell=1}^{2k-3} (\widehat{B}_{i,\ell} - \widehat{B}_{j,\ell})^2} \end{aligned}$$

- Risk when choosing Model  $M_i$ :

$$R_i = \sum_{j=1}^m \|B_i - B_j\| \cdot \Pr(M_j|D) \approx \sum_{j=1}^m \|B_i - B_j\| \cdot \frac{e^{-BIC(M_j)/2}}{\sum_{h=1}^m e^{-BIC(M_h)/2}}$$

### DT

The Decision-Theoretic (DT) criterion of Minin, Abdo, Joyce, Sullivan (2003) is to choose the model with the minimal risk

$$R_i \approx \sum_{j=1}^m \|B_i - B_j\| \cdot \frac{e^{-BIC(M_j)/2}}{\sum_{h=1}^m e^{-BIC(M_h)/2}}$$

based on the initial tree.

In a follow-up paper they study the robustness of this approach against uncertainty about the initial tree.

## Model averaging

Let  $\theta$  be the vector of parameters and  $s(\theta)$  some interesting aspect of the parameters.  $s$  must have the same meaning in all considered models  $M_1, \dots, M_m$ . We can then estimate:

$$\Pr(s(\theta)|D) \approx \frac{1}{m} \sum_{i=1}^m \Pr(s(\theta)|D, M_i) \cdot \Pr(M_i|D)$$

One possible implementation of Model averaging is reversible-jump MCMC, see Huelsenbeck, Larget, Alfaro (2004)

## Reversible-Jump MCMC

If an MCMC procedure shall sample from a state space that has several continuous components of different dimensions (e.g. for averaging over several models with different numbers of parameters), the problem arises that a density of  $n$  dimensions cannot be directly compared to a density in e.g.  $n + 1$  dimensions in a Metropolis-Hastings ratio. [1.5ex] The simple solution is to add an artificial parameter to the state of  $n$  dimensions, which has a uniform distribution on  $[0, 1]$  and no influence on the probability of the data. [1.5ex] Then you can apply Metro-Hasting to perform *reversible jumps* between the components of dimension  $n$  and dimension  $n + 1$ .

## Parametric bootstrap approach

If different models lead to different results, and it is not clear which model fits best, one should ask for all  $i$  and  $j$ :

*If model  $M_i$  was right, how accurate would an analysis based on model  $M_j$  be?*

do for each  $i$ :

1.  $\hat{\theta}_i$  := estimate  $\theta$  based on  $M_i$
2. repeat for  $k = 1, \dots, 1000$ :
  - (a)  $D_{i,k}$ : simulated data based on  $M_i$  and  $\hat{\theta}_i$
  - (b) For all  $j$ : let  $\tilde{\theta}_{i,k,j}$  the  $M_j$ -based estimation for dataset  $D_{i,k}$
3. Analyse for all  $j$  how close the average  $\overline{\tilde{\theta}_{i,..,j}}$  is to  $\hat{\theta}_i$ .

## 9.2 Does model selection matter?

### Substitution models for phylogeny reconstruction

Study with wide range of data sets for : Ripplinger and Sullivan (2008)

- Different model selection methods led to different models in 80% of the cases
- use of different best-fit models changes the optimal tree topology in 50% of the cases, but only for poorly supported branches.
- BIC and DT selected simpler models than hLRT and AIC. The simpler models performed at least as well as the more complicated.
- Use of models supported by model selection in ML gave better trees than MP or ML with K2P.
- Trees based on models favored by different model selection strategies gave similar results in hypothesis tests.
- Recommend to use the simpler BIC- and DT-selected models.

### From Lin Himmelmann's PhD thesis

Simulation study to compare (relaxed) molecular-clock models

**MC** strict molecular clock model

**CPP** compound Poisson process

**DM** Dirichlet model

**ULN** uncorrelated log-normal

**UEX** uncorrelated exponential

### ULN, UEX, DM

In ULN, UEX and DM each edge in the tree gets a rate randomly drawn from the distribution and uncorrelated to the neighboring branches.[2ex]

e.g. in the case of ULN, the logarithm of the rate on the current branch follows a normal distribution with mean  $\log(\bar{r}) + \sigma^2/2$  and variance  $\sigma^2$ , which leads to an expectation value of  $\bar{r}$  for the rate.

### The CPP model

- Rate change points are peppered randomly into the tree at rate  $\lambda$ .
- At each change point, the current rate is multiplied with  $r$ , which is drawn from a  $\Gamma$ -distribution.
- Problem: If  $\mathbb{E}r = 1$  or  $\mathbb{E}[\log r] < 0$ , rates converge to 0, and if  $\mathbb{E}[\log r] > 0$ , rates converge to  $\infty$  for long branches.
- Solution:  $\Gamma$ -parameters must lead to  $\mathbb{E}[\log r] = 0$ , and a prior on  $\lambda$  must limit the number of change-points.

### Results of Lin's model comparison

Data origin	Performance of models in analysis
MC	MC best CPP, DM, ULN almost as good UEX much worse
CPP, DM, ULN	MC, CPP, DM, ULN give good results UEX slightly worse
UEX	DM, ULN best UEX slightly worse CPP worse MC worst

Lin recommends: DM, ULN okay for most situations

## 10 Insertion-Deletion Models for Statistical Alignment

### 10.1 Alignment sampling with pairHMMs

**To Do:** estimate mutation rates from sequences

ACTCGCGCTT  
ACGTCGATT

**Classical Approach:**

1. Take **best** Alignment:

AC\_TCGCGCTT  
 ACGTCGA\_\_TT

2. Count Mutations in **best** Alignment:

1 Mismatch : 7 Matches

2 Indels (3 Sites) : 8 homologous Sites

**Problem:** underestimation of mutation rates, since alignment **fits too well!**  
 What are **typical** Alignments and Mutation Rates for given sequences?

**Idea:** Generate many random alignments  $A$  with corresponding mutation rates  $M$  according to

$$\Pr ( (A, M) \mid \text{sequences} )$$

**Needed:** A **model** of sequence evolution with insertions, deletions and substitutions. Otherwise  $\Pr(\dots)$  has no meaning!

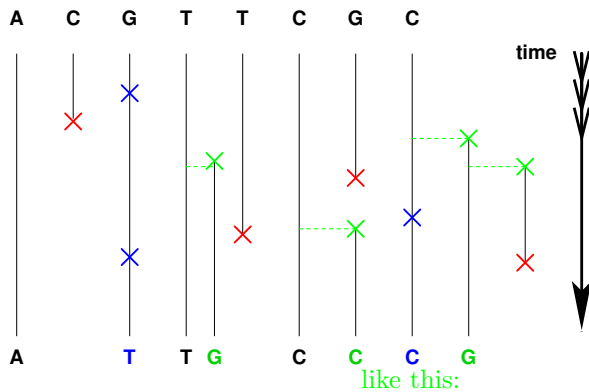
Model of Sequence Evolution

Thorne, Kishino, Felsenstein (1991):

**Deletions** with rate  $\mu$  at each site.

**Insertions** with rate  $\lambda$  right of each site & at the very left.

**Substitutions** with Rate  $s$  at each site.



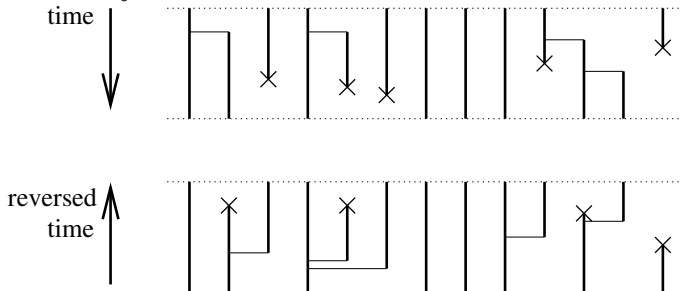
not like this:

TKF alignment convention:

ACGT\_TC\_GC\_  
 A\_TTG\_CC\_CG

ACGT\_TCG\_C\_  
 A\_TTG\_C\_CCG

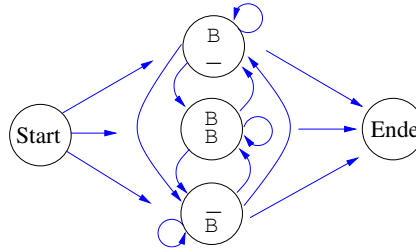
Reversibility?



### Consequence of TKF convention

The bare alignment

BBBB\_BB\_BB\_  
B\_BBB\_BB\_BB

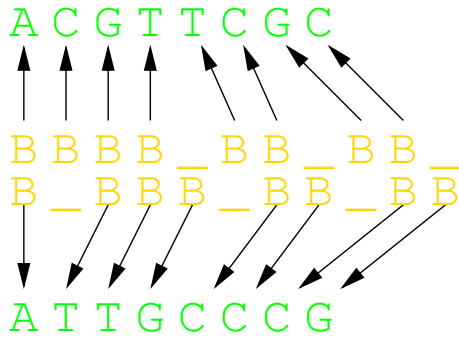


is generated by a Markov chain:

from \ to	$\bar{B}$	B	$\bar{B}$
$\bar{B}$	$(1 - \lambda\beta) \frac{\lambda}{\mu} e^{-\mu}$	$(1 - \lambda\beta) \frac{\lambda}{\mu} (1 - e^{-\mu})$	$\lambda\beta$
B	$\lambda\beta \frac{e^{-\mu}}{1 - e^{-\mu}}$	$\lambda\beta$	$\frac{1 - e^{-\mu} - \mu\beta}{1 - e^{-\mu}}$
$\bar{B}$	$(1 - \lambda\beta) \frac{\lambda}{\mu} e^{-\mu}$	$(1 - \lambda\beta) \frac{\lambda}{\mu} (1 - e^{-\mu})$	$\lambda\beta$

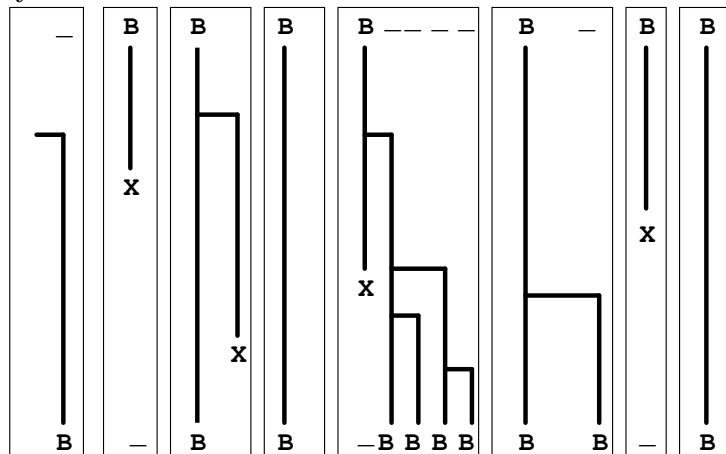
transition probabilities in (model: TKF'91),  $\beta = \frac{1 - e^{\lambda - \mu}}{\mu - \lambda e^{\lambda - \mu}}$

The Markov chain (the alignment) is hidden, observable is the pair of sequences emitted by the alignment.

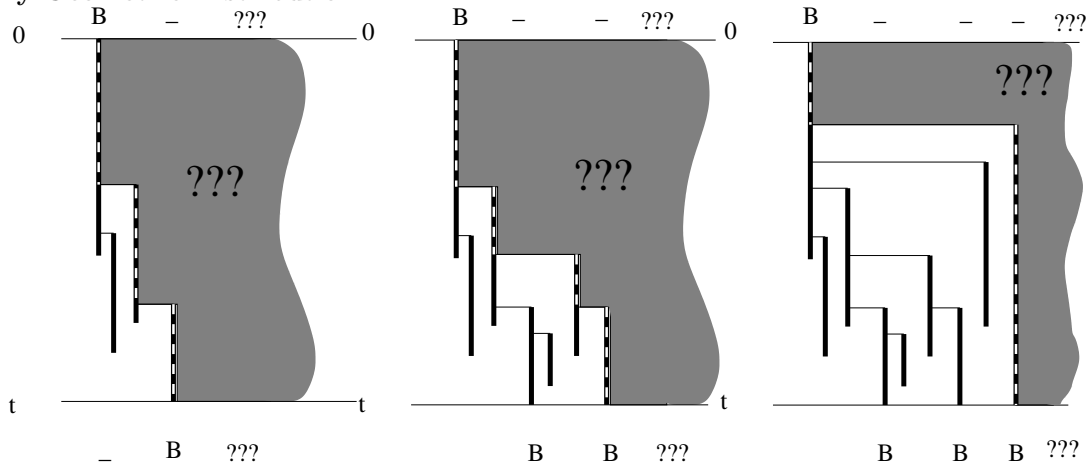


### pair Hidden Markov Model (pair HMM)

Why Markov?

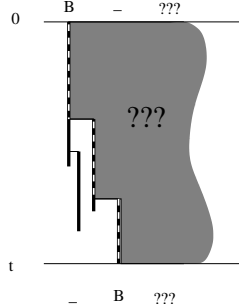


### Why Geometric Distribution?



### Computing transition probabilities

Simplification:  $\lambda = \mu$



$$\begin{aligned} \mathbb{E}(X) &= 1 \\ \Pr(X = k \mid X > 0) &= (1 - p)^{k-1} \cdot p \\ \frac{1}{p} &= \mathbb{E}(X \mid X > 0) = 1 + t \cdot \lambda \\ \Rightarrow p &= 1 / (1 + t \cdot \lambda) \end{aligned}$$

$$\begin{aligned} 1 &= \mathbb{E}(X) \\ &= \Pr(X = 0) \cdot \mathbb{E}(X \mid X = 0) \\ &\quad + \Pr(X > 0) \cdot \mathbb{E}(X \mid X > 0) \\ &= \Pr(X > 0) \cdot (1 + t \cdot \lambda) \\ \Rightarrow \Pr(X > 0) &= 1 / (1 + t \cdot \lambda) \end{aligned}$$

$X :=$  number of survivors at time  $t$

**Aim:** Sequences are given. Generate alignments  $A$  and mutation rates  $M = (\lambda, \mu, s)$  according to

$$\Pr ( (A, M) \mid \text{sequences} )$$

partial steps:

1. Assume that the mutation rates  $M$  are known. Generate alignments  $A$  according to

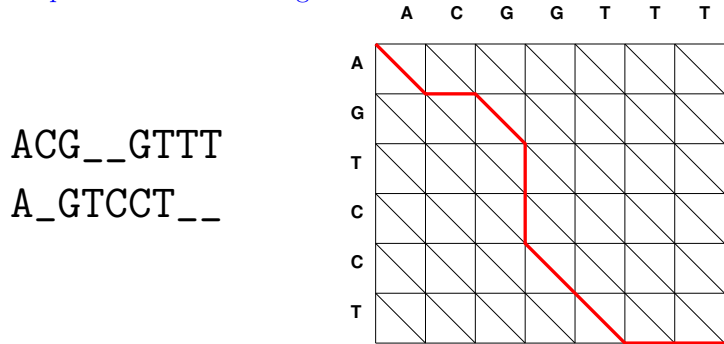
$$\Pr ( A \mid \text{sequences}, M )$$

2. Assume that the alignment  $A$  is known. Generate values for the mutation rates  $M$  according to

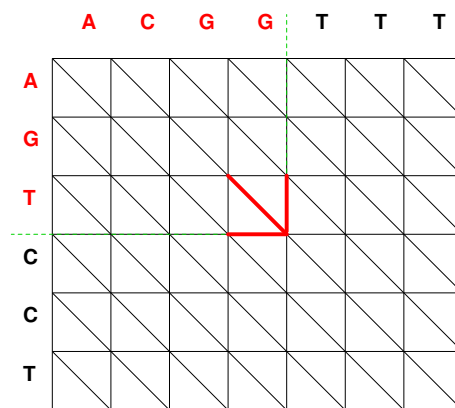
$$\Pr ( M \mid \text{sequences}, A )$$

3. combine 1. and 2.

Path representation of an alignment



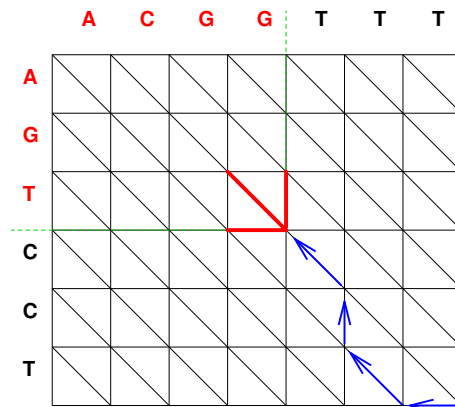
$$\Pr(\text{sequences} \mid M) = \sum_{\text{alignment } A} \Pr(A, \text{sequ.} \mid M)$$



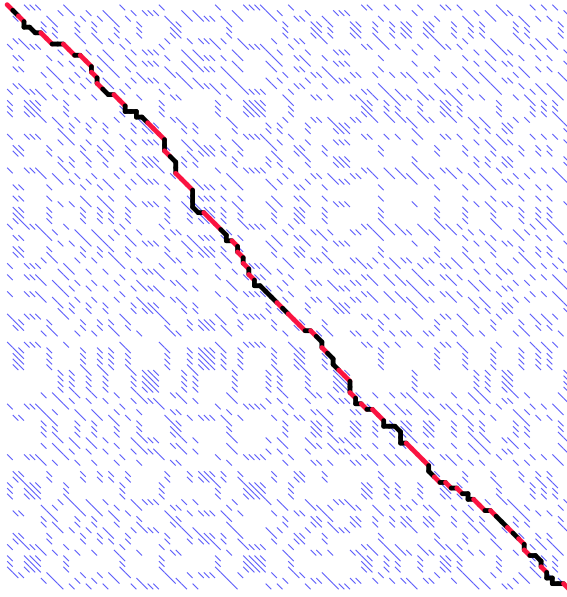
Summing efficiently: label each edge with

$\Pr(\text{Alignment contains this edge and generates the sequences so far} \mid M)$

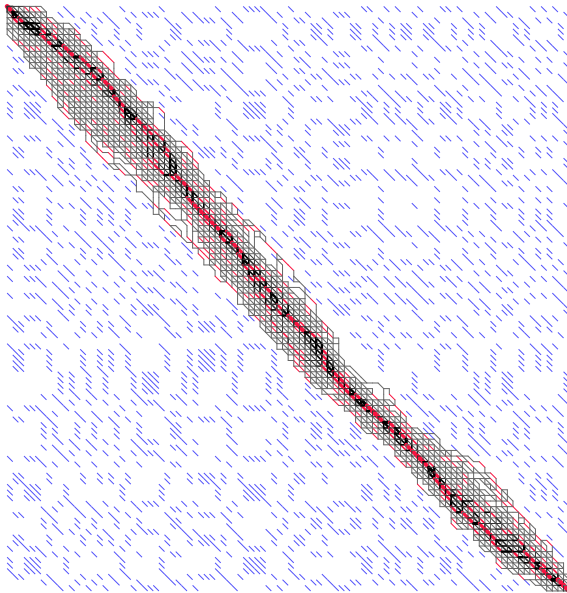
After labeling all edges, generate alignment backwards.



Random decisions in each step depend on **edge labels** and **Markov transition probabilities**.



True alignment for simulated sequence pair of length 100 with indel rate 0.3 and substitution rate 0.4.



5000 sampled alignments for simulated sequence pair of length 100 with indel rate 0.3 and substitution rate 0.4

**partial steps:**

1. Assume that the **mutation rates**  $M$  are known. Generate **alignments**  $A$  according to

$$\Pr ( A \mid \text{sequences}, M )$$

(as explained before)

2. Assume that the **alignment**  $A$  is known. Generate values for the **mutation rates**  $M$  according to

$$\Pr ( M \mid \text{sequences}, A )$$

by Metropolis Hastings Algorithm

3. Combine 1. and 2.

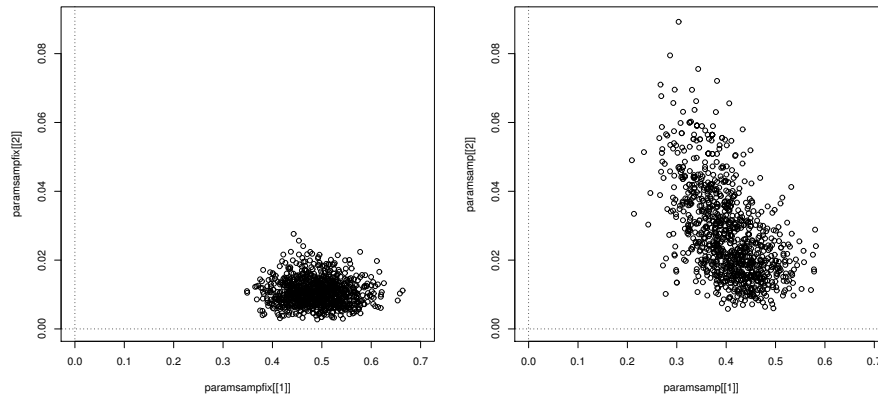
(Gibbs-Sampling)



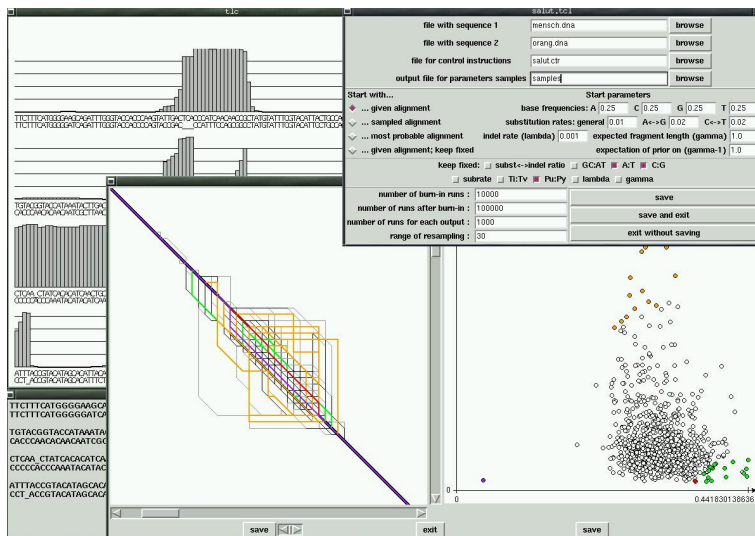
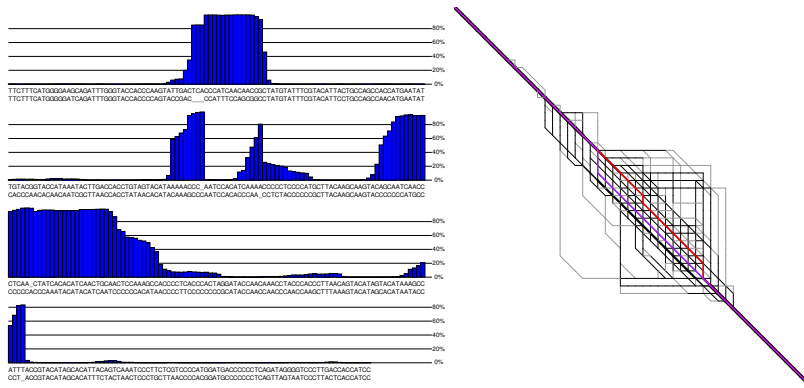
=> **Markov chain Monte Carlo Method** for sampling  $(A, M)$  according to

$$\Pr(A, M \mid \text{sequences})$$

posterior probability samplings of mutation parameters for HVR-1 of human and orangutan with alignment given in data base (left) and alignments sampled simultaneously with parameters (right)



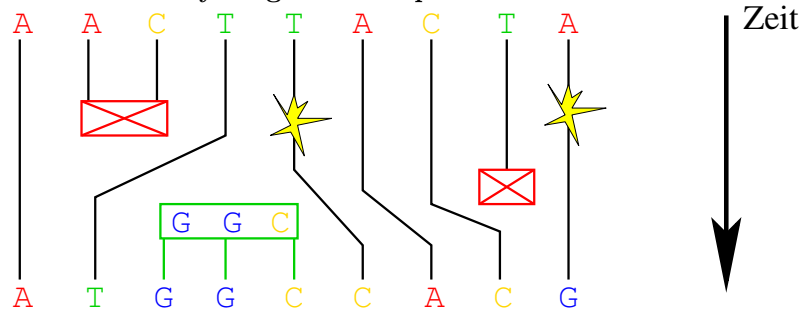
### Alignment Accuracy: HVR1 of Human and Orang



D. Metzler, R. Fleißner, A. Wakolbinger, A. von Haeseler (2001) Assessing variability by joint sampling of alignments and mutation rates, *J. Mol Evol.* 53:660-669.

## 10.2 Insertions and deletions of more than one site

InDels are usually longer than 1 position



J.L. Thorne, H. Kishino, J. Felsenstein (1992) Inching towards reality: an improved likelihood model for sequence evolution. *J. Mol. Evol.*, **34**, 3-16.

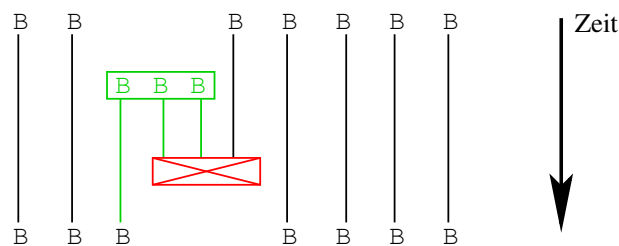
D. Metzler (2003) Statistical alignment based on fragment insertion and deletion models, *Bioinformatics* 19:490-499.

FID Model (also a pairHMM):

- instead of single nucleotides, fragments are inserted and deleted with rate  $\lambda$ .
- Length of the fragments: geometrically distributed, mean length:  $\gamma$ .

$$\Pr(L = k) = \frac{1}{\gamma} \left(1 - \frac{1}{\gamma}\right)^k$$

forbidden in TKF92 and FID:



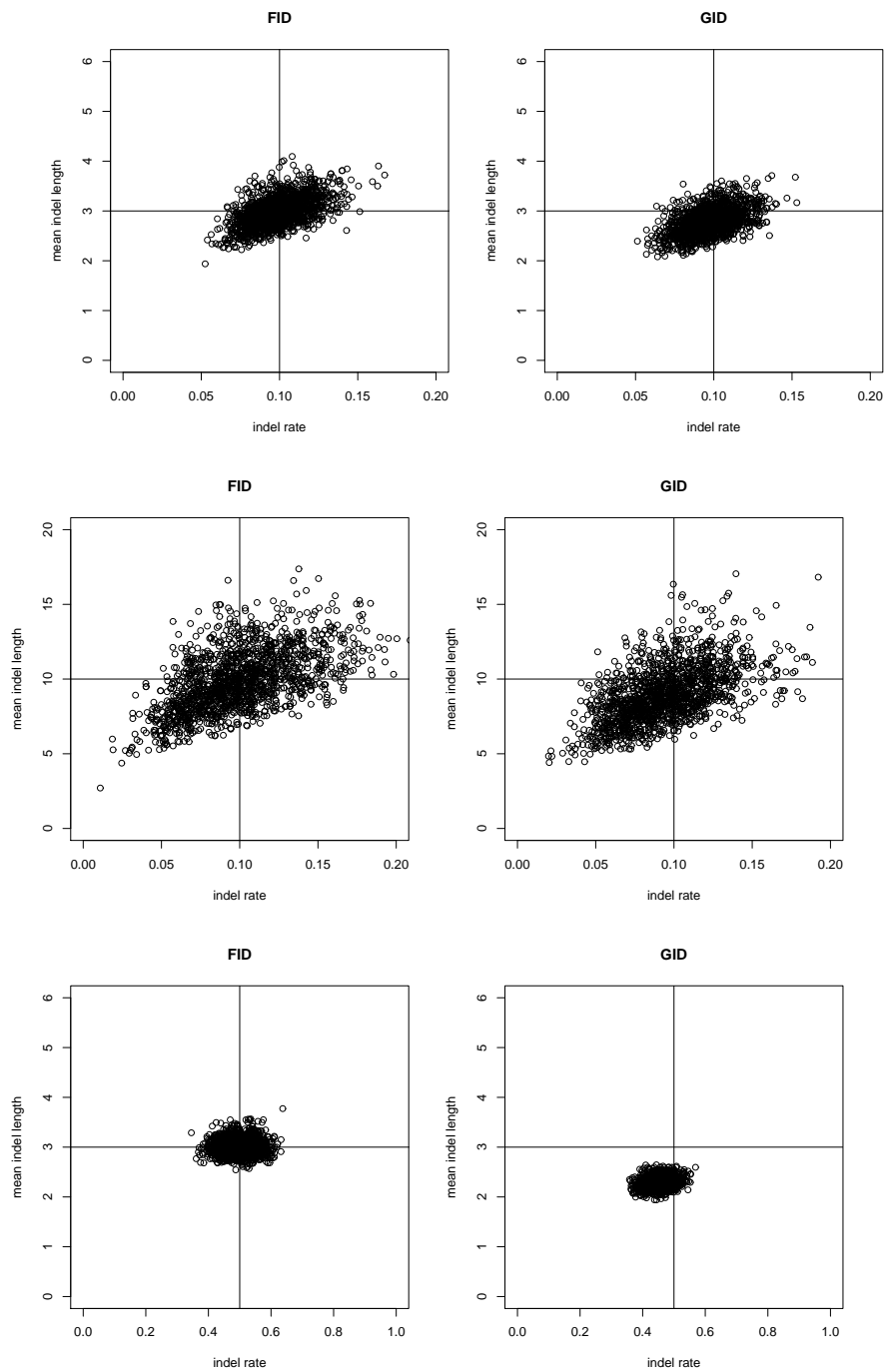
GID Model:

- ↑ this is allowed
- **no hidden Markov structure**

Use GID to simulate data and test robustness of FID

Test robustness of ML estimates for mutation rates

- Generate sequence pairs according to FID and GID
- Tell FID-based estimator which positions are homologous
- Are estimates for GID data worse than FID data? (This will be the case only when true parameter values are extreme.)
- Differences should be lower when estimates are based on sequences instead of homology structures.

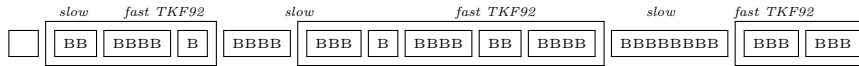


How good are FID-based methods when GID/“Long Indel Model” is true?

- no problem for parameter estimations (Metzler, 2003)
- alignment accuracy can be decreased (Miklos, Lunter, Holmes, 2004)

Maybe generate mixed-geometric gap-length with different types of fragments.  
 Along a tree fragmentation may change from edge to edge.

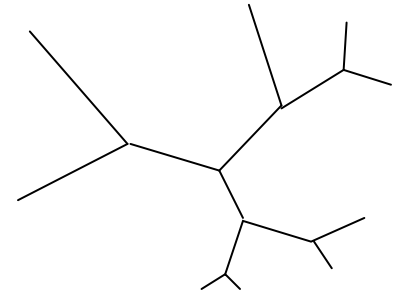
**InDel Model for detecting conserved regions**  
 A. Arribas-Gil, D. Metzler, J.-L. Plouhinec (2007)



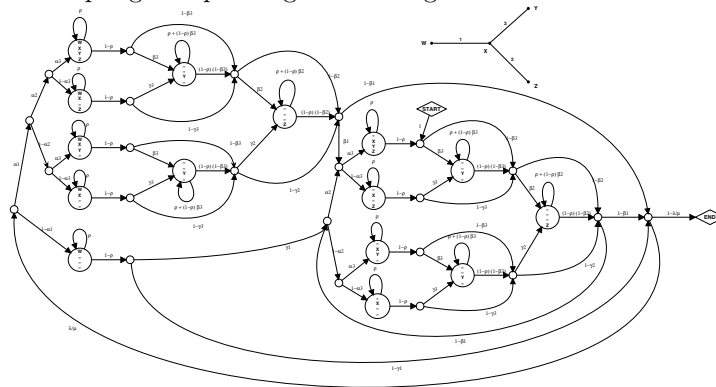
### 10.3 Multiple Alignments

I. Holmes, W. J. Bruno (2001) Evolutionary HMMs: a Bayesian approach to multiple alignment, *Bioinformatics* 17:803-820.

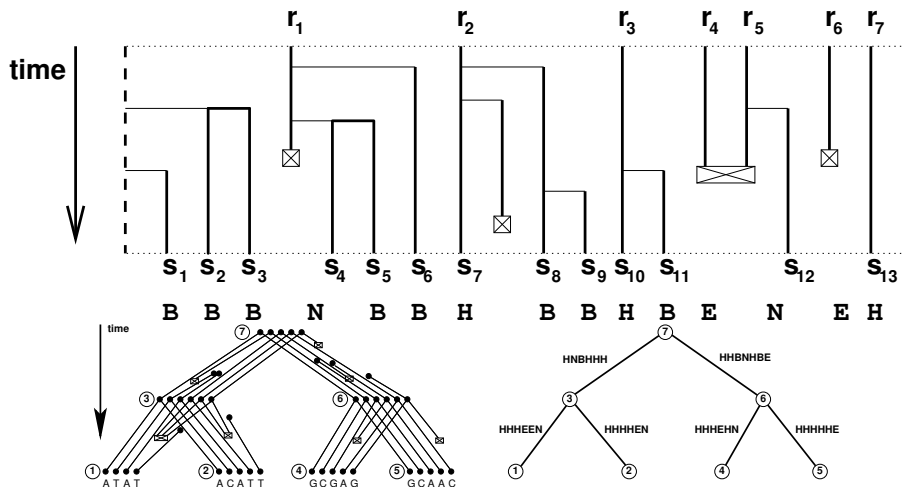
R. Fleißner, D. Metzler, A. von Haeseler (2005) Simultaneous statistical multiple alignment and phylogeny reconstruction. *Systematic Biology* 54(4):548-61.

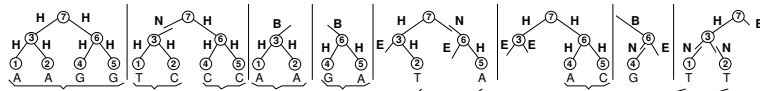


multiple HMM for sampling a sequence given its neighbours



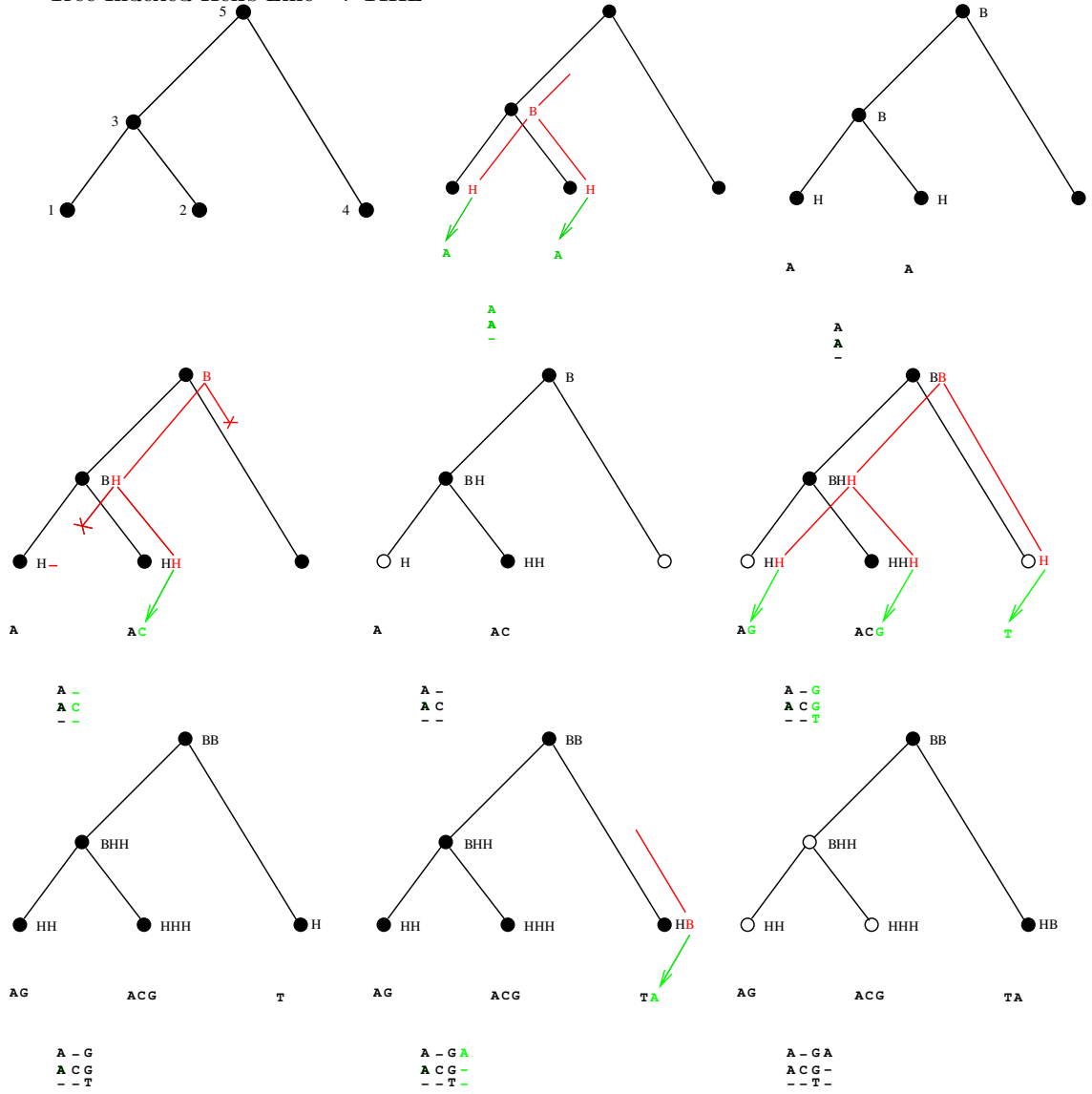
G.A. Lunter, I. Miklós, Y.S. Song, J. Hein (2003) An efficient algorithm for statistical multiple alignment on arbitrary phylogenetic trees. *J. Comp. Biol.* 10(6):869-889.

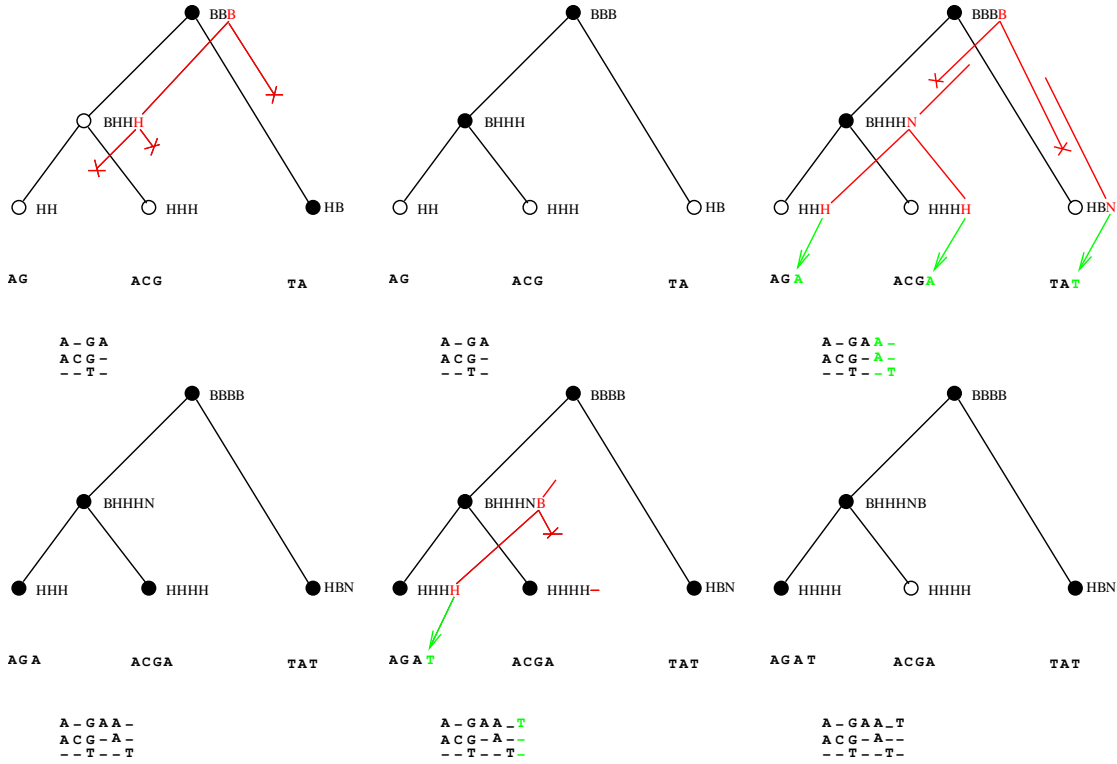




① A - T A - - - - T -  
 ② A - C A - T - - - - T  
 ④ G C - - G - A G - -  
 ⑤ G C - - A - A C - - -

Tree-Indexed Heirs Line =: TIHL





TKF91: states of hidden Markov chain are the [Sets Of Active Nodes \(soans\)](#).

$$P_S(k) = \sum_{(\mathcal{R}, e) : \mathcal{S}=[\mathcal{R}, e]} p(e)q(e)P_{\mathcal{R}}(k - v_e)\vartheta(e, k)$$

where

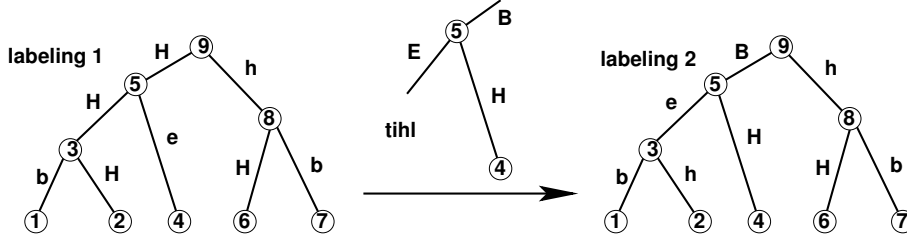
- $k$  : Multi-index of Positions in sequences at leaves
- $\mathcal{S} = [\mathcal{R}, e]$  : tihl  $e$  turns soan  $\mathcal{S}$  into soan  $\mathcal{R}$
- $P_S(k)$  : Pr(sequences up to  $k$  are generated and end there)
- $p(e)$  = Pr(indel history of  $e$ )
- $q(e)$  = Pr(no inserts at nodes in  $e$ )
- $\vartheta(e, k)$  = Pr( $e$  emits base given in data types at  $k$ )
- $v_e \in \{0, 1\}^n$  : indicates postions in leaf-sequences to which  $e$  emits

TKF91: states of hidden Markov chain are the [Sets Of Active Nodes \(soans\)](#).

**Transfer this to FID or TKF92** (fragmentation may change from edge to edge)

- D. Metzler, R. Fleißner, A. Wakolbinger, A. von Haeseler (2005) Stochastic insertion-deletion processes and statistical sequence alignment.
- D. Metzler, R. Fleißner (2007) Sequence Evolution Models for Simultaneous Alignment and Phylogeny Reconstruction.

state space: edge-labellings with  $\{B, H, e, b, h\}$ .



tihl = tree indexed heirs line

Example: 3-leaved tree

TKF91:  $2^3 = 8$  possible sets of active nodes

TKF92/FID:  $5^3 = 125$  possible labellings, 41 of them are relevant

### Why Statistical Alignment is Important

- without statistical alignment, methods like Bayesian tree sampling and bootstrapping will be by far too optimistic about the uncertainty in phylogeny inference
- over-optimization of alignments can bias your analysis
- statistical alignment is needed to use the information contributed by insertions and deletions

## 10.4 Software for joint estimation of phylogenies and alignments

### BAlI-Phy

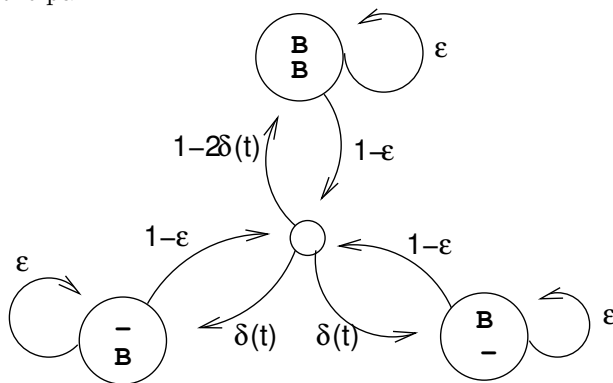
<http://www.biomath.ucla.edu/msuchard/bali-phy/>

## References

- [RS05] B.D. Redelings, M.A. Suchard (2005) Joint Bayesian Estimation of Alignment and Phylogeny *Systematic Biology* **54(3)**:401-418
- [SR06] M.A. Suchard, B.D. Redelings (2006) BAlI-Phy: simultaneous Bayesian inference of alignment and phylogeny *Bioinformatics* **22**:2047-2048
- [RS07] B.D. Redelings, M.A. Suchard (2007) Incorporating indel information into phylogeny estimation for rapidly emerging pathogens. *BMC Evolutionary Biology* **7**:40

### pairHMM in BAlI-Phy

The alignment consists of a geometrically distributed number of fragments. It is generated according to the pairHMM



with  
 $\delta(t) = 1 - e^{-\lambda t / (1-\epsilon)}$ .

## MCMC steps in BAli-Phy

- Parts of the pairwise alignments along branches of the current tree are re-sampled. Felsenstein wild-cards are used for the nucleotide or amino acid types, i.e. probability distributions conditioned on the sequences at the tips of the tree.
- SPR steps for updating the tree.
- After an SPR step a pairwise alignment along the new branch is sampled. For efficiency, it keeps the alignments within each of the two partial trees fixed.

## Other software for statistical alignment

**StatAlign** provides a graphical user interface where you can watch the changes in the alignment and the phylogeny

<http://phylogeny-cafe.elte.hu/StatAlign/>

**Alifritz** does not sample from the posterior distribution but searches the alignment and the phylogeny of the highest posterior probability

<http://www.cibiv.at/software/alifritz/>

## References

- [NMLH08] A. Novák, I. Miklós, R. Lyngsø, J. Hein (2008) StatAlign: an extendable software package for joint Bayesian estimation of alignments and evolutionary trees. *Bioinformatics* **24(20)**: 2403-2404
- [FMvH] R. Fleißner, D. Metzler, A. von Haeseler (2005) Simultaneous Statistical Multiple Alignment and Phylogeny Reconstruction. *Syst. Biol.* **54**: 548-561

# 11 Quantitative Characters and Independent Contrasts

## 11.1 Brownian motions along the branches of the tree

### Type of questions to be answered

Quantitative traits like number of genes, mutation rates, or morphological traits like weight or body length differ for different species.

- Do two traits evolve in a correlated way or are their values just correlated because they evolved independently along the same tree?
- Is a trait significantly different for a certain trait such that adaptation must have played a role?
- Can we use morphological traits for phylogeny reconstruction?

Model for the neutral evolution of a quantitative trait along the branches of a phylogenetic tree.

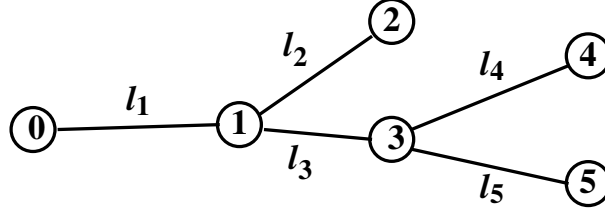
- Independent on different branches
- After an appropriate rescaling it changes randomly like a Brownian motion. This is a Markov process with

$$X_{s+t} - X_s \sim \mathcal{N}(0, t),$$

where  $\mathcal{N}(0, t)$  means normal distribution with mean 0 and variance  $\sigma^2 = t$ .



Example: Brownian motion starts in node 0 of this tree with a non-random value  $X_0$ :



Then,  $\mathbb{E}X_i = x_0$  for all  $i$ , and the variance of any node is its distance to the root, e.g.  $\text{var}(X_5) = l_1 + l_3 + l_5$ .

$$\begin{aligned} \text{cov}(X_5, X_4) &= \text{cov}(X_5 - X_3 + X_3, X_4 - X_3 + X_3) \\ &= \text{cov}(X_5 - X_3, X_4 - X_3) + \text{cov}(X_5 - X_3, X_3) \\ &\quad + \text{cov}(X_3, X_4 - X_3) + \text{cov}(X_3, X_3) \\ &= \text{var}(X_3) = l_1 + l_3 \end{aligned}$$

In general: The covariance of the values  $X_k$  and  $X_\ell$  at the nodes  $k$  and  $\ell$  is the variance  $\text{var}(X_h)$  of the value at their most recent common ancestor  $h$ .

Let  $v_i$  be the parent node of node  $i$ , then the values  $\left(\frac{X_i - X_{v_i}}{\sqrt{\ell_i}}\right)_{i=1, \dots, n}$  are stochastically independent and standard-normally distributed. Together they are a standard-normally distributed random vector.

Moreover, the map

$$Y := \begin{pmatrix} \frac{X_1 - X_{v_1}}{\sqrt{\ell_1}} \\ \vdots \\ \frac{X_n - X_{v_n}}{\sqrt{\ell_n}} \end{pmatrix} \mapsto \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix} = X$$

is an affine transformation, i.e. can be represented as  $Y \mapsto w + MY = X$  with appropriate vector  $w$  and matrix  $M$ . This implies that  $X$  is also normally distributed, and its distribution is determined by its expectation value and its covariance matrix.

## 11.2 Excursus: Multidimensional Normal Distribution

- The vector (height,width) is a two-dimensional random vector (with values in  $\mathbb{R}^2$ )
- An  $d$ -dimensional random vector is a vector of  $d$  random elements
- The expectation of a random vector  $X = (X_1, X_2, \dots, X_d)^T$  is the vector of the expectations:

$$\mathbb{E}X = \mathbb{E} \begin{pmatrix} X_1 \\ \vdots \\ X_d \end{pmatrix} = \begin{pmatrix} \mathbb{E}X_1 \\ \vdots \\ \mathbb{E}X_d \end{pmatrix}$$

- The expectation of a random matrix  $M = (M_{ij})_{i=1..n, j=1..d}$  is the matrix of the expectations:

$$\mathbb{E} \begin{pmatrix} M_{11} & M_{12} & \cdots & M_{1d} \\ \vdots & \vdots & \ddots & \vdots \\ M_{n1} & M_{n2} & \cdots & M_{nd} \end{pmatrix} = \begin{pmatrix} \mathbb{E}M_{11} & \mathbb{E}M_{12} & \cdots & \mathbb{E}M_{1d} \\ \vdots & \ddots & & \vdots \\ \mathbb{E}M_{n1} & \mathbb{E}M_{n2} & \cdots & \mathbb{E}M_{nd} \end{pmatrix}$$

- reminder: The variance of a univariate random variable  $X$  is  $\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}X)^2] = \mathbb{E}[X^2] - (\mathbb{E}X)^2$ .
- The analog in the multivariate case is the so called *covariance matrix* (or dispersion matrix or variance-

covariance matrix). The covariance matrix  $\text{Var}(X) = \Sigma$  of  $X = (X_1, \dots, X_d)^T$  is

$$\begin{aligned}\Sigma &= \begin{pmatrix} \text{Cov}(X_1, X_1) & \text{Cov}(X_1, X_2) & \cdots & \text{Cov}(X_1, X_d) \\ \text{Cov}(X_2, X_1) & \text{Cov}(X_2, X_2) & \cdots & \text{Cov}(X_2, X_d) \\ \vdots & & \ddots & \\ \text{Cov}(X_d, X_1) & \text{Cov}(X_d, X_2) & \cdots & \text{Cov}(X_d, X_d) \end{pmatrix} \\ &= \mathbb{E} \left[ \begin{pmatrix} X_1 - \mathbb{E}X_1 \\ \vdots \\ X_d - \mathbb{E}X_d \end{pmatrix} \cdot (X_1 - \mathbb{E}X_1, \dots, X_d - \mathbb{E}X_d) \right] \\ &= \mathbb{E} \left[ (X - \mathbb{E}X) \cdot (X - \mathbb{E}X)^T \right] \\ &= \mathbb{E} \left[ X \cdot X^T \right] - \mathbb{E}X \cdot (\mathbb{E}X)^T\end{aligned}$$

- Linearity of the expectation is analogous to the univariate case: Let  $X = (X_1, \dots, X_d)$  be a random vector and  $C = (C_{ij})_{i=1..n, j=1..d}$  be a deterministic matrix. Then

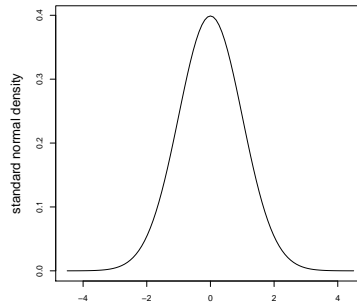
$$\mathbb{E}(C \cdot X) = C \cdot \mathbb{E}(X)$$

- If  $Y := X - \mathbb{E}(X)$ , then

$$\begin{aligned}\text{Var}(C \cdot X) &= \text{Var}(C \cdot Y) \\ &= \mathbb{E} [C \cdot Y \cdot (C \cdot Y)^T] \\ &= \mathbb{E} [C \cdot Y \cdot Y^T \cdot C^T] \\ &= C \cdot \mathbb{E} [Y \cdot Y^T] \cdot C^T \\ &= C \cdot \text{Var}(Y) \cdot C^T \\ &= C \cdot \text{Var}(X) \cdot C^T\end{aligned}$$

- Reminder: Univariate normal distribution  $\mathcal{N}(\mu, \sigma^2)$  with mean  $\mu \in \mathbb{R}$  and variance  $\sigma^2 \in (0, \infty)$  has the density

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$



Remember:  $\Pr(\mu - \sigma < X < \mu + \sigma) = 0.68$  and  $\Pr(\mu - 1.96\sigma < X < \mu + 1.96\sigma) = 0.95$

- The density of the  $d$ -dimensional normal distribution with mean  $\mu \in \mathbb{R}^d$  and covariance matrix  $\Sigma \in \mathbb{R}^{d \times d}$  is analogous:

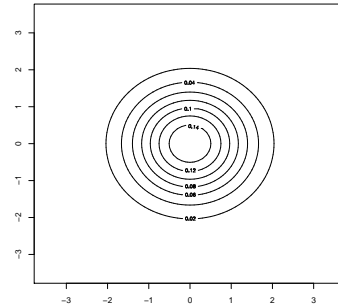
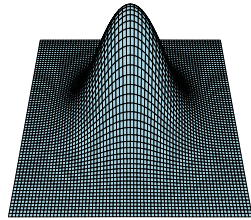
$$f(x) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp\left(-\frac{(x-\mu)^T \Sigma^{-1} (x-\mu)}{2}\right)$$

for  $x \in \mathbb{R}^d$  where  $\det(\Sigma)$  is the determinant of  $\Sigma$ , and  $\Sigma^{-1}$  is the inverse matrix. We write  $\mathcal{N}_d(\mu, \Sigma)$  for this distribution.

- The *standard multivariate normal distribution* has mean  $\mu = 0$  and the identity matrix  $\Sigma = \mathbb{I}$  as covariance matrix.

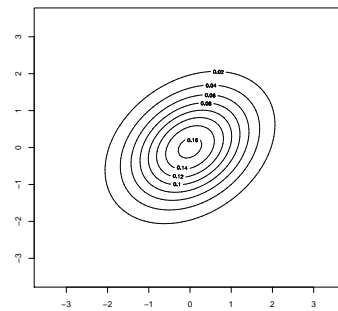
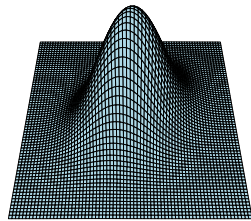
**Plots for  $d = 2$**

Correlation 0.0:  $\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ ,  $\text{Var}(X_1) = 1 = \text{Var}(X_2)$ ,  $\text{Cov}(X_1, X_2) = 0.0$



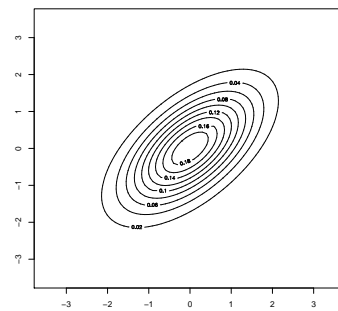
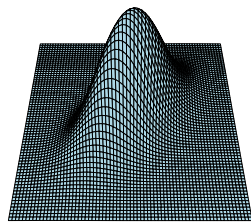
**Plots for  $d = 2$**

Correlation 0.3:  $\Sigma = \begin{pmatrix} 1 & 0.3 \\ 0.3 & 1 \end{pmatrix}$ ,  $\text{Var}(X_1) = 1 = \text{Var}(X_2)$ ,  $\text{Cov}(X_1, X_2) = 0.3$



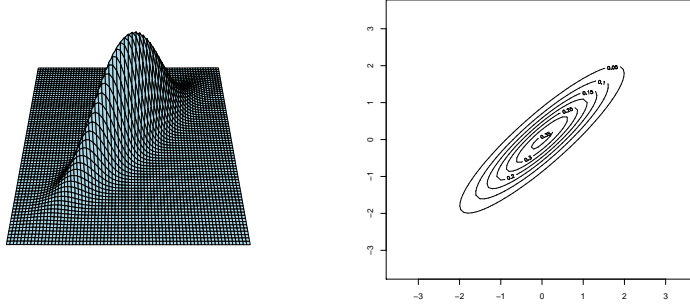
**Plots for  $d = 2$**

Correlation 0.6:  $\Sigma = \begin{pmatrix} 1 & 0.6 \\ 0.6 & 1 \end{pmatrix}$ ,  $\text{Var}(X_1) = 1 = \text{Var}(X_2)$ ,  $\text{Cov}(X_1, X_2) = 0.6$



**Plots for  $d = 2$**

Correlation 0.9:  $\Sigma = \begin{pmatrix} 1 & 0.9 \\ 0.9 & 1 \end{pmatrix}$ ,  $\text{Var}(X_1) = 1 = \text{Var}(X_2)$ ,  $\text{Cov}(X_1, X_2) = 0.9$



Properties: (Let the distribution of  $X$  be  $\mathcal{N}_d(\mu, \Sigma)$ )

- Linear combinations are univariate normal distributed:  $\langle c, X \rangle \sim \mathcal{N}(\langle c, \mu \rangle, c \Sigma c^T)$
- $X_i$  and  $X_j$  are independent  $\iff \text{Cov}(X_i, X_j) = 0$
- The standardized normal distribution is standard normal distributed

$$\Sigma^{-\frac{1}{2}} \cdot (X - \mu) \sim \mathcal{N}_d(0, \mathbb{I})$$

where  $M = \Sigma^{-\frac{1}{2}}$  is a matrix such that  $M^T \cdot M \cdot \Sigma = \mathbb{I}$ .

- The square of the standardized normal distribution is chi-squared distributed with  $d$  degrees of freedom:

$$(X - \mu)^T \Sigma^{-1} (X - \mu) \sim \chi_d^2.$$

- If  $Y_1, Y_2, \dots, Y_d$  are independent and standard normal distributed, then  $(Y_1, \dots, Y_d) \sim \mathcal{N}(0, \mathbb{I})$ .

Properties: (Let the distribution of  $X$  be  $\mathcal{N}_d(\mu, \Sigma)$ )

- If  $M \in \mathbb{R}^{p \times d}$  is a non-random matrix, then  $M \cdot X \sim \mathcal{N}_p(M \cdot \mu, M \Sigma M^T)$

### 11.3 Why to use REML

Assume now that the values of  $X_i$  in the tips of the tree are given and that the topology of the tree is known. How can we estimate the branch lengths? Let's apply ML!

Example: For a rooted tree with two tips, we measure the values  $x_{1i}$  and  $x_{2i}$  for  $i = 1, \dots, p$  of  $p$  different traits in the tips 1 and 2. The values  $x_{0i}$  in the root of the tree are unknown. For known values  $\sigma_i$  we assume that the value of trait  $x_{ji}$  for  $j \in \{1, 2\}$  is normally distributed with mean  $x_{0i}$  and variance  $\ell_j \sigma_i^2$ , where  $\ell_j$  is the unknown length of the branch to tip  $j$ . We have to maximize the likelihood

$$\begin{aligned} L(x_0, \ell_1, \ell_2) &= \prod_{i=1}^p \frac{1}{\sqrt{2\pi\sigma_i^2\ell_1}} \cdot e^{-\frac{(x_{1i}-x_{0i})^2}{2\ell_1\sigma_i^2}} \cdot \frac{1}{\sqrt{2\pi\sigma_i^2\ell_2}} \cdot e^{-\frac{(x_{2i}-x_{0i})^2}{2\ell_2\sigma_i^2}} \\ &= \prod_{i=1}^p \frac{1}{2\pi\sigma_i^2\sqrt{\ell_1\ell_2}} \cdot e^{-\frac{1}{2\sigma_i^2} \left( \frac{(x_{1i}-x_{0i})^2}{\ell_1} + \frac{(x_{2i}-x_{0i})^2}{\ell_2} \right)} \\ &= \frac{1}{\prod_{i=1}^p \sigma_i^2} \cdot \left( \frac{1}{2\pi\sqrt{\ell_1\ell_2}} \right)^p \cdot e^{-\frac{1}{2} \cdot \left( \sum_{i=1}^p \frac{1}{\sigma_i^2} \cdot \left( \frac{(x_{1i}-x_{0i})^2}{\ell_1} + \frac{(x_{2i}-x_{0i})^2}{\ell_2} \right) \right)} \end{aligned}$$

To find values  $x_{01}, \dots, x_{0p}$  and  $\ell_1$  and  $\ell_2$  that maximize  $L(x_0, \ell_1, \ell_2)$ , we first minimize the exponents

$$\frac{(x_{1i} - x_{0i})^2}{\ell_1} + \frac{(x_{2i} - x_{0i})^2}{\ell_2}$$

by setting

$$\widehat{x_{0i}} = \frac{\frac{x_{1i}}{\ell_1} + \frac{x_{2i}}{\ell_2}}{\frac{1}{\ell_1} + \frac{1}{\ell_2}}$$

Then we search for  $\ell_1$  and  $\ell_2$  that maximize

$$\frac{(x_{1i} - x_{0i})^2}{\ell_1} + \frac{(x_{2i} - x_{0i})^2}{\ell_2}$$

This means that  $\ell_1\ell_2$  should be small and  $\ell_1 + \ell_2$  should be large, and we get that by setting  $\ell_1 = 0$  and  $\ell_2 = \infty$  or vice versa.

This is perhaps not what we expected. What is the reason for this absurd result?

Heuristic explanation: We have one parameter per  $i$  too much in the model. This parameter vanishes for  $\ell_1 = 0$  because this forces  $x_{0i}$  to be  $x_{1i}$ .

There are several ways to circumvent this problem, which also appears when for trees with more than two tips, e.g.:

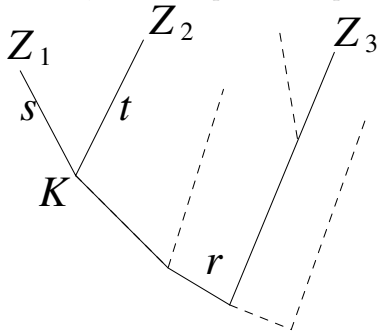
- Assume a strict molecular clock such that all tips must have the same distance to the root (Thompson, 1975)
- Felsenstein's REML (REduced Maximum-Likelihood) approach is to avoid the root and consider only unrooted trees. For the example above this means that we only estimate  $\ell_1 + \ell_2$  by the ML estimator

$$\widehat{\ell_1 + \ell_2} = \frac{1}{p} \sum_{i=1}^p \left( \frac{x_{1i} - x_{2i}}{\sigma_i} \right)^2.$$

## 11.4 Computing Independent Contrasts by Pruning the Tree

Let  $Z = (Z_1, \dots, Z_m)^T$  be the vector of values for a quantitative character in the tips  $b_1, \dots, b_m$  of the tree. To compute the likelihood of the tree or correct correlations for phylogenetic relationship or to decide whether there is significant evidence for adaptation, we apply REML and transform the values in the tips back into a standard-normally distributed vector.

One way of doing this is a variant of Felsenstein's pruning algorithm. It leads to *independent* transformations – so-called *contrasts* – between the values in the tips that can be associated with the branches of the tree, which helps to interpret them.



We start with the contrast  $Z_2 - Z_1$ . Then we assign a value  $W$  to node  $k$  (the MRCA of nodes  $b_1$  and  $b_2$ ) that is a weighted average of  $Z_1$  and  $Z_2$  but independent of the contrast  $Z_2 - Z_1$ : Set

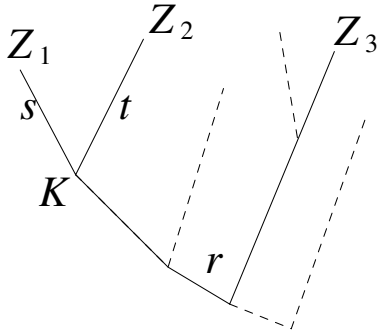
$$W := x \cdot Z_1 + (1 - x) \cdot Z_2$$

and search for  $x$  such that

$$0 = \text{cov}(x \cdot Z_1 + (1 - x) \cdot Z_2, Z_1 - Z_2)$$

$$\begin{aligned} &= x \cdot \text{var}(Z_1) - x \cdot \text{cov}(Z_1, Z_2) + (1 - x) \cdot \text{cov}(Z_2, Z_1) - \\ &\quad (1 - x) \cdot \text{var}(Z_2) \\ &= x \cdot \text{var}(Z_1) - x \cdot \text{var}(K) + (1 - x) \cdot \text{var}(K) - (1 - x) \cdot \text{var}(Z_2) \\ &= x \cdot \text{var}(Z_1 - K) - (1 - x) \cdot \text{var}(Z_2 - K) \\ &= x \cdot s - (1 - x) \cdot t \end{aligned}$$

$$\Rightarrow x = \frac{t}{s + t}$$



Hence, we set

$$W := \frac{t}{s+t} \cdot Z_1 + \frac{s}{s+t} \cdot Z_2.$$

If the distance between  $k$  and some tip with value  $Z_3$  is  $r$ , then

$$\text{var}(K - Z_3) = r,$$

where  $K$  is the value in node  $k$ .

should not consider  $W$  as an estimate for  $K$  because  $\text{var}(W - Z_3)$

$$\begin{aligned} &= \text{var} \left( \frac{s}{t+s} \cdot Z_1 + \frac{t}{t+s} \cdot Z_2 - Z_3 \right) \\ &= \text{var} \left( \frac{s}{t+s} \cdot (Z_1 - K) + \frac{t}{t+s} \cdot (Z_2 - K) + K - Z_3 \right) \\ &= \left( \frac{s}{s+t} \right)^2 \cdot \text{var}(Z_1 - K) + \left( \frac{t}{s+t} \right)^2 \cdot \text{var}(Z_2 - K) + \text{var}(K - Z_3) \\ &= \frac{s^2}{(s+t)^2} \cdot t + \frac{t^2}{(s+t)^2} \cdot s + r = \frac{st}{s+t} + r > \text{var}(K - Z_3). \end{aligned}$$

Thus, we can imagine that we prune the subtree of  $b_1$  and  $b_2$  from the tree and extending the branch to  $k$  by length  $\frac{st}{s+t}$ . To the new tip at the end of this extended branch we assign a value of  $W$ . The contrast  $Z_1 - Z_2$  is uncorrelated to all values at tips of this new tree and thus also to any contrasts that we can compute from them.

We repeat this pruning step until we have  $m - 1$  independent contrasts. Dividing all contrasts by their standard deviations leads to a standard-normally distributed vector of contrasts.

All this is only true under the null hypothesis of neutral evolution. We can reject this null hypothesis if the vector of standard-normalized contrast deviates significantly from the normal distribution. Since the contrasts are associated with branches of the tree, we can then identify which branch of the tree shows evidence for process of adaptation. (Here we assume that the phylogeny is known.)

In principle, we can also use quantitative characters to estimate the tree, but usually the amount of available data is insufficient to infer the tree, adaptation processes and correlation between different quantitative traits. It usually makes more sense to estimate the tree from molecular data and then use the independent contrasts method to analyse the evolution of the quantitative traits along the tree.

## 11.5 Software

### Phylip: contrasts

<http://evolution.genetics.washington.edu/phylip/doc/contrast.html>

Can deal with variation of traits within species (Above we have always assumed only one value for per trait for each species. This should be the average value, which, however, can usually not be estimated with high precision.)

Note that correlation of different traits within species is usually different from correlation between species.

## References

- [F08] J. Felsenstein (2008) Comparative Methods with Sampling Error and Within-Species Variation: Contrasts Revisited and Revised. *American Naturalist* **171**(6): 713–725

**BayesTraits**

The software package BayesTraits from Mark Pagel's group provides several Bayesian and Likelihood-based methods for inferring the evolution of continuous and discrete traits along phylogenetic trees.

<http://www.evolution.reading.ac.uk/BayesTraits.html>