

Estimating the Relative Order of Speciation or Coalescence Events on a Given Phylogeny

Tanja Gernhard¹, Daniel Ford², Rutger Vos³ and Mike Steel⁴

¹Department of Mathematics, Kombinatorische Geometrie (M9), TU München, Boltzmannstr. 3, 85747 Garching, Germany.

²Department of Mathematics, Stanford University, U.S.A.

³Department of Biological Sciences, Simon Fraser University, Vancouver, Canada.

⁴Biomathematics Research Centre, University of Canterbury, Christchurch, New Zealand.

Abstract: The reconstruction of large phylogenetic trees from data that violates clocklike evolution (or as a supertree constructed from any m input trees) raises a difficult question for biologists—how can one assign relative dates to the vertices of the tree? In this paper we investigate this problem, assuming a uniform distribution on the order of the inner vertices of the tree (which includes, but is more general than, the popular Yule distribution on trees). We derive fast algorithms for computing the probability that (i) any given vertex in the tree was the j -th speciation event (for each j), and (ii) any one given vertex is earlier in the tree than a second given vertex. We show how the first algorithm can be used to calculate the expected length of any given interior edge in any given tree that has been generated under either a constant-rate speciation model, or the coalescent model.

Keywords: Phylogenetics, neutral model, dating speciation events, edge lengths.

1. Introduction

A fundamental task in evolutionary biology is constructing evolutionary trees from a variety of data. These constructed trees show the ancestral relationship between the species.

Not only the relationship between species is of interest, but also the time between speciation events. When constructing an evolutionary tree from a set of molecular data which satisfies the molecular clock, the edge lengths can be interpreted as a time scale. In many cases, no time scale is obtained when constructing a tree though:

- Often, molecular data does not satisfy the molecular clock and so the edge lengths do not represent a time scale.
- Trees can be constructed from morphological data or non-standard molecular data like gene order. This does not provide any edge lengths.
- Having several different trees, one can combine them and construct a ‘supertree.’ Even though there may have been time scales on the original trees, most supertree methods return a tree without a time scale.

For those trees, we still want to find edge lengths representing the time between speciation events. In this paper, we will estimate the edge lengths from the shape of the tree. The method works for trees which evolved under the Yule model [Yule, 1924; Edwards, 1970; Harding, 1971; Page, 1991]. Under the Yule model, in each point of time, each species is equally likely to split. Minor changes to the method for the Yule model give us an edge length estimation for trees under the popular coalescent setting [Nordborg, 2001].

An example for a tree with unknown edge lengths is the primate supertree \mathcal{T}_p recently published in [Vos and Mooers]. Figure 1 shows a part of \mathcal{T}_p . The primate tree is a supertree on 218 species and was constructed with the MRP method (Matrix Representation using Parsimony analysis, see [Baum, 1992; Ragan, 1992]). Since for most of the interior vertices, no molecular estimates were available, the edge lengths for the tree were estimated. In [Vos and Mooers], 10^6 rank functions on \mathcal{T}_p were drawn uniformly at random. For each of those rank functions, the expected time intervals, i.e. the edge lengths, between vertices were considered (the expected waiting time after the $(n - 1)$ th event until the n th event is $1/n$).

Correspondence: Tanja Gernhard, Department of Mathematics, Kombinatorische Geometrie (M9), TU München, Boltzmannstr. 3, 85747 Garching, Germany. Tel: +49 89 289 16882; Email: gernhard@ma.tum.de

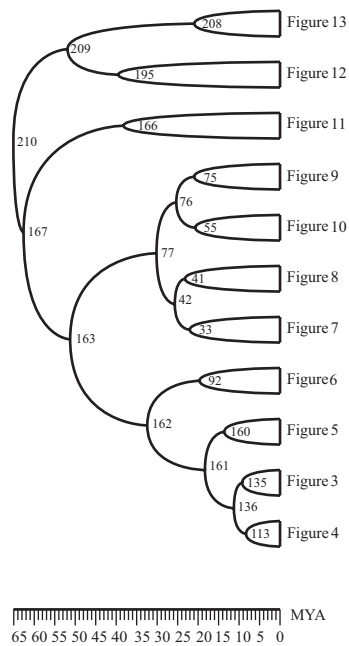


Figure 1. Part of the primate supertree. Figure 4–13 are some subtrees, for details see [Vos and Mooers].

The authors of [Vos and Mooers] concluded their paper by asking for an analytical approach to the estimation of the edge length, which we will provide below.

In order to estimate the edge lengths, we developed the algorithms RANKPROB and COMPARE. Those algorithms answer questions like:

Was speciation event with label 76 in the primate tree (see Fig. 1) more likely to be an early event in the tree or a late event? What is the probability that 76 was the 6th speciation event? Was it more likely that speciation event 76 happened before speciation event 162 or 162 before 76?

The algorithms work for trees where every labeled history is equiprobable. This class of model, which includes the Yule model and the coalescent model, has been popular in macroevolutionary studies [Nee and May, 1997; Zhaxybayeva and Gogarten, 2004]. Note that the algorithms here are the same for the Yule model and the coalescent model, whereas the edge length estimation has minor differences for the two models.

The algorithms RANKPROB, COMPARE and an algorithm for obtaining the expected rank and variance for a vertex were implemented in Python, see [Gernhard, 2006].

2. Probability Distribution of the Rank of a Vertex

Let \mathcal{T} be a rooted phylogenetic tree [Semple and Steel, 2003] with $|V| = n$ leaves. The set of interior vertices of \mathcal{T} shall be $\overset{\circ}{V}$. For a binary tree, we have $|\overset{\circ}{V}| = n - 1$. Let the function r be a bijection from the set of interior vertices $\overset{\circ}{V}$ of \mathcal{T} into $\{1, 2, \dots, |\overset{\circ}{V}|\}$ with $r(v_1) \leq r(v_2)$ if v_1 is an ancestor of v_2 . The function r is called a *rank function* for \mathcal{T} . A vertex v with $r(v) = i$ is said to have *rank* i . Note that r induces a linear order on the set $\overset{\circ}{V}$. Further, define $r(\mathcal{T}) := \{r : r \text{ is a rank function on } \mathcal{T}\}$. We are interested in the distribution of the possible ranks for a certain vertex, i.e. we want to know the probability of $r(v) = i$ for a given $v \in \overset{\circ}{V}$. If every rank function on a given tree is equally likely, we have

$$\mathbb{P}[r(v) = i] = \frac{|\{r : r(v) = i, r \in r(\mathcal{T})\}|}{|r(\mathcal{T})|} \quad (1)$$

which will be calculated for rooted binary trees in polynomial time by algorithm RANKPROB. In the algorithm, we will use the formula [Semple and Steel, 2003]

$$|r(\mathcal{T})| = \frac{|\overset{\circ}{V}|!}{\prod_{v \in \overset{\circ}{V}} (n_v - 1)} \quad (2)$$

where n_v is the number of leaves below v . Note that Equation 2 holds for binary and nonbinary trees.

Examples of stochastic models on phylogenetic trees where each rank function is equally likely include:

- The Yule model has the probability distribution $\mathbb{P}[r|\mathcal{T}] = \frac{\prod_{v \in \overset{\circ}{V}} (n_v - 1)}{(n-1)!}$ which is the uniform distribution [Edwards, 1970; Brown, 1994].
- The coalescent model has the same probability distribution on rooted binary ranked trees as the Yule model. So $\mathbb{P}[r|\mathcal{T}]$ is the uniform distribution [Aldous, 2001].
- For some sets of trees (e.g. those drawn from the uniform model [Pinelis, 2003], also known as PDA model), no rank function is induced. If one assumes that all rank functions are equally likely on these trees, one can apply Equation 1 to such trees as well.

2.1. A polynomial-time algorithm

The following algorithm calculates the probability distribution of the rank of a vertex v in a rooted

binary phylogenetic tree \mathcal{T} . The idea of the algorithm is the following (cf. Figure 2). Label the vertices on the path from v to the root ρ by $v = x_1, \dots, x_n = \rho$. Let \mathcal{T}_m be the subtree of \mathcal{T} containing the vertex x_m and all its descendants. Let $\alpha_{\mathcal{T}_m, v}(i)$ be the number of rank functions on the tree \mathcal{T}_m where v has rank i . The values $\alpha_{\mathcal{T}_m, v}(i)$, $i = 1, \dots, |\mathring{V}|$ are calculated iteratively for $m = 1, \dots, n$. The probability $\mathbb{P}[r(v) = i]$ equals $\frac{\alpha_{\mathcal{T}_n, v}(i)}{\sum_{i=1}^{|\mathring{V}|} \alpha_{\mathcal{T}_n, v}(i)}$. The α -values in the fraction have a lot of factors in common which cancel out. In the following algorithm, we calculate α -values without the unnecessary terms instead, $\tilde{\alpha}_{\mathcal{T}_m, v}(i)$. We have $\alpha_{\mathcal{T}_m, v}(i) = \tilde{\alpha}_{\mathcal{T}_m, v}(i) |r(\mathcal{T}_1)| |r(\mathcal{T}'_1)| \dots |r(\mathcal{T}'_{m-1})|$.

Algorithm: RANKPROB(\mathcal{T}, v)

Input: A rooted binary phylogenetic tree \mathcal{T} and an interior vertex v .

Output: The probabilities $\mathbb{P}[r(v) = i]$ for $i = 1, \dots, |\mathring{V}|$.

- 1: Denote the vertices of the path from v to root ρ with $(v = x_1, x_2, \dots, x_n = \rho)$.
- 2: Denote the subtree of \mathcal{T} , consisting of root x_m and all its descendants, by \mathcal{T}_m for $m = 1, \dots, n$.
- 3: Initialize $\tilde{\alpha}_{\mathcal{T}_m, v}(i) := 0$ for $i = 1, \dots, |\mathring{V}_{\mathcal{T}}|$, $m = 1, \dots, n$.
- 4: $\tilde{\alpha}_{\mathcal{T}_1, v}(1) := 1$
- 5: **for** $m = 2, \dots, n$ **do**
- 6: $\mathcal{T}'_{m-1} := \mathcal{T}_m \setminus (\mathcal{T}_{m-1} \cup x_m)$ (cf. Figure 3)

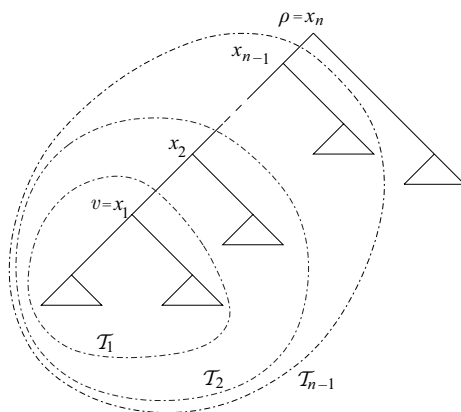


Figure 2. Labeling the tree for the algorithm RANKPROB.

- 7: **for** $i = m, \dots, |\mathring{V}_{\mathcal{T}_m}|$ **do**
- 8: $M := \min \{ |\mathring{V}_{\mathcal{T}'_{m-1}}|, i - 2 \}$
- 9: $\tilde{\alpha}_{\mathcal{T}_m, v}(i) := \sum_{j=0}^M \tilde{\alpha}_{\mathcal{T}_{m-1}, v}(i - j - 1) \binom{|\mathring{V}_{\mathcal{T}_{m-1}}| + |\mathring{V}_{\mathcal{T}'_{m-1}}| - (i - 1)}{|\mathring{V}_{\mathcal{T}'_{m-1}}| - j} \binom{i - 2}{j} (*)$
- 10: **end for**
- 11: **end for**
- 12: **for** $i = 1, \dots, |\mathring{V}_{\mathcal{T}}|$ **do**
- 13: $\mathbb{P}[r(v) = i] := \frac{\tilde{\alpha}_{\mathcal{T}_n, v}(i)}{\sum_j \tilde{\alpha}_{\mathcal{T}_n, v}(j)}$
- 14: **end for**
- 15: RETURN $\mathbb{P}[r(v) = i]$, $i = 1, \dots, |\mathring{V}|$.

Proving the correctness and runtime of RANKPROB makes use of the following two observations.

Remark 1. Let A_i be a set containing n_i elements with a linear order, $i \in \{1, 2\}$. There are $\binom{n_1+n_2}{n_1}$ possible linear orders on $A_1 \cup A_2$ which preserve the linear order on A_1 and A_2 . This follows from the observation that the number of such linear orders on $A_1 \cup A_2$ is equivalent to the number of ways of choosing n_1 elements from $n_1 + n_2$ elements, which is $\binom{n_1+n_2}{n_1}$.

Remark 2. The values $\binom{n}{k}$ for all $n, k \leq N$ ($n, k, N \in \mathbb{N}$) can be calculated in $O(N^2)$ using Pascal's Triangle. Thus, after $O(N^2)$ calculations, any value $\binom{n}{k}$ with $n, k \leq N$ can be obtained in constant time.

Theorem 3. RANKPROB returns the quantities

$$\mathbb{P}[r(v) = i]$$

for each given $v \in \mathring{V}$ and all $i \in 1, \dots, |\mathring{V}|$. The runtime is $O(|\mathring{V}|^2)$.

Proof. Let $\alpha_{\mathcal{T}_m, v}(i) = \tilde{\alpha}_{\mathcal{T}_m, v}(i) |r(\mathcal{T}_1)| |r(\mathcal{T}'_1)| |r(\mathcal{T}'_2)| \dots |r(\mathcal{T}'_{m-1})|$. We first show that $\alpha_{\mathcal{T}_m, v}(i) = |\{r : r(v) = i, r \in r(\mathcal{T}_m)\}|$ for $m = 1, \dots, n$, $i = 1, \dots, |\mathring{V}_{\mathcal{T}}|$. That implies

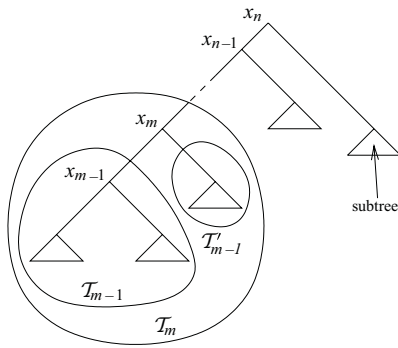


Figure 3. Labeling the tree for the recursion in RANKPROB.

$$\begin{aligned} \mathbb{P}[r(v) = i] &= \frac{|\{r : r(v) = i, r \in r(\mathbf{T})\}|}{|r(\mathbf{T})|} \\ &= \frac{\alpha_{\mathbf{T},v}(i)}{\sum_i \alpha_{\mathbf{T},v}(i)} = \frac{\tilde{\alpha}_{\mathbf{T},v}(i)}{\sum_i \tilde{\alpha}_{\mathbf{T},v}(i)} \end{aligned}$$

which proves the theorem.

The proof is by induction over m .

For $m = 1$, $\alpha_{\mathcal{T}_1,v}(1) = |\{r : r(v) = 1, r \in r(\mathcal{T})\}|$. Vertex v is the root of \mathcal{T}_1 , so $\alpha_{\mathcal{T}_1,v}(i) = 0$ for all $i > 1$.

Let $m = k$ and $\alpha_{\mathcal{T}_m,v}(i) = |\{r : r(v) = i, r \in r(\mathcal{T}_m)\}|$ holds for all $m < k$. $\alpha_{\mathcal{T}_k,v}(i) = 0$ clearly holds for all $i > |\mathring{V}_{\mathcal{T}_k}|$ since $r_{\mathcal{T}_k} : v \rightarrow \{1, \dots, |\mathring{V}_{\mathcal{T}_k}|\}$. So it remains to verify that the term (*) returns the right values for $\alpha_{\mathcal{T}_k,v}(i)$. Assume that the vertex v is in the $(i - j - 1)$ -th position in \mathcal{T}_{k-1} (with $i - j - 1 > 0$) for some rank function $r_{\mathcal{T}_{k-1}}$ and v shall be in the i -th position in \mathcal{T}_k .

Now combine the linear order in the tree \mathcal{T}_{k-1} induced by $r_{\mathcal{T}_{k-1}}$ with a linear order in \mathcal{T}'_{k-1} induced by $r_{\mathcal{T}'_{k-1}}$ to get a linear order on \mathcal{T}_k . The first j vertices of \mathcal{T}'_{k-1} must be inserted between vertices of \mathcal{T}_{k-1} with lower rank than v so that v ends up to be in the i -th position of the tree \mathcal{T}_k . Count the number of possible way to do this as follows. The tree \mathcal{T}'_{k-1} has $|r(\mathcal{T}'_{k-1})|$ possible rank functions. Combining a rank function $r_{\mathcal{T}_{k-1}}$ with a rank function $r_{\mathcal{T}'_{k-1}}$ to get a rank function $r_{\mathcal{T}_k}$ with $r_{\mathcal{T}_k}(v) = i$ means inserting the first j vertices of \mathcal{T}'_{k-1} anywhere between the first $(i - j - 2)$ vertices of \mathcal{T}_{k-1} . There are

$$\binom{(i - j - 2) + j}{j} = \binom{i - 2}{j}$$

possibilities according to Remark 1. For combining the $|\mathring{V}_{\mathcal{T}_{k-1}}| - (i - j - 1)$ vertices of rank bigger than

v in \mathcal{T}_{k-1} with the remaining $|\mathring{V}_{\mathcal{T}'_{k-1}}| - j$ vertices in \mathcal{T}'_{k-1} , there are

$$\begin{aligned} &\binom{|\mathring{V}_{\mathcal{T}_{k-1}}| - (i - j - 1) + |\mathring{V}_{\mathcal{T}'_{k-1}}| - j}{|\mathring{V}_{\mathcal{T}'_{k-1}}| - j} \\ &= \binom{|\mathring{V}_{\mathcal{T}_{k-1}}| + |\mathring{V}_{\mathcal{T}'_{k-1}}| - (i - 1)}{|\mathring{V}_{\mathcal{T}'_{k-1}}| - j} \end{aligned}$$

possibilities. This follows again from Remark 1. The number of rank functions $r_{\mathcal{T}_{k-1}}$ with $r_{\mathcal{T}_{k-1}}(v) = i - j - 1$ is $\alpha_{\mathcal{T}_{k-1},v}(i - j - 1)$ by the induction assumption. Multiplying all those possibilities gives

$$\begin{aligned} &\alpha_{\mathcal{T}_{k-1},v}(i - j - 1) |r(\mathcal{T}'_{k-1})| \\ &\binom{|\mathring{V}_{\mathcal{T}_{k-1}}| + |\mathring{V}_{\mathcal{T}'_{k-1}}| - (i - 1)}{|\mathring{V}_{\mathcal{T}'_{k-1}}| - j} \binom{i - 2}{j} \end{aligned}$$

where $\alpha_{\mathcal{T}_{k-1},v}(i) = \tilde{\alpha}_{\mathcal{T}_{k-1},v}(i) |r(\mathcal{T}_1)| |r(\mathcal{T}'_1)| |r(\mathcal{T}'_2)| \dots |r(\mathcal{T}'_{k-2})|$. The value $|\{r : r(v) = i, r \in r(\mathcal{T})\}|$ is then the sum over all possible j which establishes the correctness of the algorithm.

All that remains is to verify the runtime. Note that the combinatorial factors $\binom{n}{k}$ for all $n, k \leq |V|$ can be calculated in advance in quadratic time, see Remark 2. In the algorithm, those factors can then be obtained in constant time.

The most time consuming part of the algorithm is line 13. Adding up all calculations needed for obtaining $\alpha'_{\mathcal{T}_m,v}(i)$, $m = 1, \dots, n$, $i = 1, \dots, |\mathring{V}_{\mathcal{T}_m}|$ comes to:

$$\begin{aligned} \sum_{m=2}^n |\mathring{V}_{\mathcal{T}_m}| |\mathring{V}_{\mathcal{T}'_{m-1}}| &\leq \sum_{m=2}^n |\mathring{V}| |\mathring{V}_{\mathcal{T}'_{m-1}}| \\ &= |\mathring{V}| \sum_{m=2}^n |\mathring{V}_{\mathcal{T}'_{m-1}}| \leq |\mathring{V}|^2 \end{aligned}$$

The last inequality holds since the vertices of the \mathcal{T}'_m , $m = 1, \dots, n - 1$, are distinct. Therefore, the runtime is quadratic.

Remark 4. With $\mathbb{P}[r(v) = i]$ from Theorem 3, the expected value $\mu_{r(v)}$ and the variance $\sigma_{r(v)}^2$ for $r(v)$ can be calculated by

$$\mu_{r(v)} = \sum_{i=1}^{|\mathring{V}|} i \mathbb{P}[r(v) = i]$$

$$\sigma_{r(v)}^2 = \sum_{i=1}^{|\mathring{V}|} i^2 \mathbb{P}[r(v) = i] - \mu_{r(v)}^2$$

Remark 5. The algorithm RankProb can be generalized to non-binary trees [Gernhard, 2006]. The runtime is again quadratic.

3. Application of RANKPROB - Estimating Edge Lengths

3.1. The Yule model

A very common stochastic model for rooted binary phylogenetic trees with edge lengths is the continuous-time Yule model [Edwards, 1970]. As in the discrete Yule model, at every point in time, each species is equally likely to split and give birth to two new species. The expected waiting time for the next speciation event in a tree with n leaves is $1/n$. That is, each species at any given time has a constant speciation rate (normalized so that 1 is the expected time until it next speciates).

Assume that the primate tree \mathcal{T}_p evolved under the continuous-time Yule model. In [Gernhard, 2006], the tree shape of \mathcal{T}_p (i.e. the tree without edge lengths) under the discrete Yule model is tested against the uniform model and accepts the Yule model.

Here, we describe how to estimate the edge lengths for a tree which is assumed to have evolved under the continuous-time Yule model.

Let (u, v) be an interior edge in \mathcal{T} with u the immediate ancestor of v . Let X be the random variable ‘length of the edge (u, v) ’ given that \mathcal{T} is generated according to the continuous-time Yule model.

The expected length $\mathbb{E}[X]$ of the edge (u, v) is given by

$$\mathbb{E}[X] = \sum_{i,j} \mathbb{E}[X|r(u) = i, r(v) = j] \mathbb{P}[r(u) = i, r(v) = j].$$

Since, under the continuous-time Yule model, the expected waiting time for the next speciation event is $1/n$ it follows that:

$$\mathbb{E}[X|r(u) = i, r(v) = j] = \sum_{k=1}^{j-i} \frac{1}{i+k}.$$

It remains to calculate the probability $\mathbb{P}[r(u) = i, r(v) = j]$. This is equivalent to counting all the possible rank functions where $r(u) = i$ and $r(v) = j$. The subtree \mathcal{T}_v consists of v and all its descendants. The tree \mathcal{T}_u equals the tree \mathcal{T} where all the descendants of v are deleted, i.e. v is a leaf in \mathcal{T}_u , see Figure 4.

Note that $\mathbb{P}[r(u) = i, r(v) = j] = 0$ if $|\mathring{V}_{\mathcal{T}_u}| < j - 1$. Therefore, assume $|\mathring{V}_{\mathcal{T}_u}| \geq j - 1$ in the following.

The number of rank functions on \mathcal{T}_u is $|r(\mathcal{T}_u)|$. The probability $\mathbb{P}[r(u) = i]$ can be calculated with $\text{RANKPROB}(\mathcal{T}_u, u)$. So the number of rank functions in \mathcal{T}_u with $\mathbb{P}[r(u) = i]$ is $\mathbb{P}[r(u) = i] \cdot |r(\mathcal{T}_u)|$.

The number of rank functions on \mathcal{T}_v is $|r(\mathcal{T}_v)|$. Let any linear order on the trees \mathcal{T}_u and \mathcal{T}_v be given. Combining those two linear orders into an order, r , on \mathcal{T} with $r(v) = j$ means that the vertices with rank $1, 2, \dots, j - 1$ in \mathcal{T}_u keep their rank. Vertex v gets rank j . The remaining $|\mathring{V}_{\mathcal{T}_u}| - (j - 1)$ vertices in \mathcal{T}_u and $|\mathring{V}_{\mathcal{T}_v}| - 1$ vertices in \mathcal{T}_v have to be shuffled together. According to Remark (1), this can be done in

$$\binom{|\mathring{V}_{\mathcal{T}_u}| - (j - 1) + |\mathring{V}_{\mathcal{T}_v}| - 1}{|\mathring{V}_{\mathcal{T}_v}| - 1} = \binom{|\mathring{V}_{\mathcal{T}_u}| + |\mathring{V}_{\mathcal{T}_v}| - j}{|\mathring{V}_{\mathcal{T}_v}| - 1}$$

different ways. Thus overall there are:

$$\mathbb{P}[r(u) = i] \cdot |r(\mathcal{T}_u)| \cdot |r(\mathcal{T}_v)| \cdot \binom{|\mathring{V}_{\mathcal{T}_u}| + |\mathring{V}_{\mathcal{T}_v}| - j}{|\mathring{V}_{\mathcal{T}_v}| - 1}$$

different rank functions on \mathcal{T} with $r(u) = i$ and $r(v) = j$. For the probability $\mathbb{P}[r(u) = i, r(v) = j]$:

$$\frac{\mathbb{P}[r(u) = i, r(v) = j] \cdot |r(\mathcal{T}_u)| \cdot |r(\mathcal{T}_v)| \cdot \binom{|\mathring{V}_{\mathcal{T}_u}| + |\mathring{V}_{\mathcal{T}_v}| - j}{|\mathring{V}_{\mathcal{T}_v}| - 1}}{\sum_{i,j} \mathbb{P}[r(u) = i, r(v) = j] \cdot |r(\mathcal{T}_u)| \cdot |r(\mathcal{T}_v)| \cdot \binom{|\mathring{V}_{\mathcal{T}_u}| + |\mathring{V}_{\mathcal{T}_v}| - j}{|\mathring{V}_{\mathcal{T}_v}| - 1}}$$

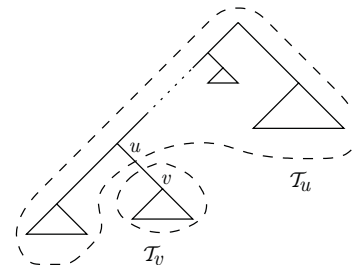


Figure 4. Labeling the tree for estimating the edge lengths.

Since $|r(\mathcal{T}_u)|$ and $|r(\mathcal{T}_v)|$ are independent of i and j , those factors cancel out, giving

$$\mathbb{P}[r(u) = i, r(v) = j] = \frac{\mathbb{P}[r(u) = i] \cdot \binom{|\dot{V}_{\mathcal{T}_u}| + |\dot{V}_{\mathcal{T}_v}| - j}{|\dot{V}_{\mathcal{T}_v}| - 1}}{\sum_{i,j} \mathbb{P}[r(u) = i] \cdot \binom{|\dot{V}_{\mathcal{T}_u}| + |\dot{V}_{\mathcal{T}_v}| - j}{|\dot{V}_{\mathcal{T}_v}| - 1}} \quad (3)$$

Furthermore, note that

$$\binom{|\dot{V}_{\mathcal{T}_u}| + |\dot{V}_{\mathcal{T}_v}| - j}{|\dot{V}_{\mathcal{T}_v}| - 1} = \frac{(|\dot{V}_{\mathcal{T}_v}| - j)!}{(|\dot{V}_{\mathcal{T}_v}| - 1)! (|\dot{V}_{\mathcal{T}}| - j - (|\dot{V}_{\mathcal{T}_v}| - 1))!}$$

Again, since $(|\dot{V}_{\mathcal{T}_v}| - 1)!$ is independent of i and j , this factor cancels out, and so

$$\mathbb{P}[r(u) = i, r(v) = j] = \frac{\mathbb{P}[r(u) = i] \cdot \prod_{k=0}^{|\dot{V}_{\mathcal{T}_v}|-2} (|\dot{V}_{\mathcal{T}}| - j - k)}{\sum_{i,j} \mathbb{P}[r(u) = i] \cdot \prod_{k=0}^{|\dot{V}_{\mathcal{T}_v}|-2} (|\dot{V}_{\mathcal{T}}| - j - k)}$$

Let $\Omega = \{(i, j) : i < j, i, j \in \{1, \dots, |\dot{V}|\}, |\dot{V}_{\mathcal{T}_u}| \geq j - 1\}$. With this notation, the expected edge length $\mathbb{E}[X]$ is

$$\begin{aligned} \mathbb{E}[X] &= \sum_{(i,j) \in \Omega} \mathbb{E}[X | r(u) = i, r(v) = j] \mathbb{P}[r(u) = i, r(v) = j] \\ &= \sum_{(i,j) \in \Omega} \left[\left(\sum_{k=1}^{j-i} \frac{1}{i+k} \right) \frac{\mathbb{P}[r(u) = i] \cdot \prod_{k=0}^{|\dot{V}_{\mathcal{T}_v}|-2} (|\dot{V}_{\mathcal{T}}| - j - k)}{\sum_{(i,j) \in \Omega} \left[\mathbb{P}[r(u) = i] \cdot \prod_{k=0}^{|\dot{V}_{\mathcal{T}_v}|-2} (|\dot{V}_{\mathcal{T}}| - j - k) \right]} \right] \\ &= \frac{\sum_{(i,j) \in \Omega} \left[\left(\sum_{k=1}^{j-i} \frac{1}{i+k} \right) \cdot \mathbb{P}[r(u) = i] \cdot \prod_{k=0}^{|\dot{V}_{\mathcal{T}_v}|-2} (|\dot{V}_{\mathcal{T}}| - j - k) \right]}{\sum_{(i,j) \in \Omega} \left[\mathbb{P}[r(u) = i] \cdot \prod_{k=0}^{|\dot{V}_{\mathcal{T}_v}|-2} (|\dot{V}_{\mathcal{T}}| - j - k) \right]} \quad (4) \end{aligned}$$

Remark 6. Equation 4 enables the estimation of the length of every interior edge. For pendant edges, the approach above gives no definite answer. All we know is that the time from the latest interior vertex, which has rank $n - 1$, until today is expected to be at most $1/n$ where n is the number of leaves.

Suppose that the growth process is stopped as soon as the $n - 1$ -st speciation event occurs. In this case the expected length X of a pendant edge below an interior vertex v is:

$$\mathbb{E}[X] = \sum_{i=1}^{n-1} \mathbb{P}[r(v) = i] \sum_{k=i}^{n-2} \frac{1}{k+1}$$

The expected depth of vertex v from the first branchpoint is:

$$\sum_{i=1}^{n-1} \mathbb{P}[r(v) = i] \sum_{k=1}^{i-1} \frac{1}{k+1}$$

So the depth Y of the leaf in question from the first branchpoint has expectation independent of v :

$$\begin{aligned} \mathbb{E}[Y] &= \sum_{i=1}^{n-1} \mathbb{P}[r(v) = i] \sum_{k=1}^{i-1} \frac{1}{k+1} \\ &+ \sum_{i=1}^{n-1} \mathbb{P}[r(v) = i] \sum_{k=i}^{n-2} \frac{1}{k+1} \\ &= \sum_{i=1}^{n-1} \mathbb{P}[r(v) = i] \sum_{k=1}^{n-2} \frac{1}{k+1} \\ &= \sum_{k=1}^{n-2} \frac{1}{k+1} \end{aligned}$$

In other words, assigning to each edge of a given tree topology its expected length gives a tree which obeys a molecular clock.

Remark 7. Often, an inferred tree has vertices with more than two descendants, i.e. there is lack of resolution due to, e.g. conflicting data. Our calculation for the expected edge length assumes a binary tree though.

However, the expected edge length may be calculated for each possible binary resolution of the supertree. Assume the supertree \mathcal{T} has the possible binary resolutions $\mathcal{T}_1, \dots, \mathcal{T}_m$. For an edge (u, v) in \mathcal{T} where u is the immediate ancestor of v , the expected edge length is calculated in the trees \mathcal{T}_i for $i = 1, \dots, m$. The expected edge length in \mathcal{T}_i is denoted by e_i for $i = 1, \dots, m$. Note that if u is a vertex with more than two descendants in \mathcal{T} then v is in general not a direct descendant of u in \mathcal{T}_i . The value e_i in resolution \mathcal{T}_i is then the sum of all expected edge lengths on the path from u to v in \mathcal{T}_i .

Calculate the expected edge length $\mathbb{E}[X]$ of (u, v) in the supertree \mathcal{T} by

$$\mathbb{E}[X] = \frac{\sum_i e_i \mathbb{P}[\mathcal{T}_i]}{\sum_i \mathbb{P}[\mathcal{T}_i]} \quad (5)$$

where the probability of a tree \mathcal{T} under the Yule model is [Brown, 1994]

$$\mathbb{P}[\mathcal{T}] = \frac{2^{n-1}}{n! \prod_{v \in \mathcal{V}} (n_v - 1)}$$

Again, once the expected length of pendant edges is included the resulting tree obeys a molecular clock, meaning that all leaves are at the same depth.

3.2. The coalescent process

The edge length estimation in the previous section works for the continuous-time Yule model. By changing the method above slightly, we get an edge length estimation for the coalescent process. In the coalescent setting, we have

$$\mathbb{E}[X | r(u)=i, r(v)=j] = \sum_{k=1}^{j-i} \frac{1}{(i+k)(i+k-1)}.$$

Therefore, the expected edge length for an interior edge (u, v) can be calculated by the following modification of Equation 4:

$$\mathbb{E}[X] = \frac{\sum_{(i,j) \in \Omega} \left[\left(\sum_{k=1}^{j-i} \frac{1}{(i+k)(i+k-1)} \right) \cdot \mathbb{P}[r(u)=i] \cdot \Pi_{k=0}^{|\dot{V}_{\mathcal{T}_v}|-2} (|\dot{V}_{\mathcal{T}}| - j - k) \right]}{\sum_{(i,j) \in \Omega} \left[\mathbb{P}[r(u)=i] \cdot \Pi_{k=0}^{|\dot{V}_{\mathcal{T}_v}|-2} (|\dot{V}_{\mathcal{T}}| - j - k) \right]}$$

The calculations in Section 3.1 and 3.2 provide exact values for the expected length of an interior edge under the Yule or coalescent process as an alternative to simulations. However simulations also provide some indication of the variability in the estimate of edge lengths, and it may be of interest to also investigate analytically the variance (or even the distribution) of the edge length in future work, rather than just its mean.

4. Comparing Two Interior Vertices

The algorithm RANKPROB can also be used for comparing two interior vertices. Assume again that every rank function on a rooted binary phylogenetic tree \mathcal{T} is equally likely. The aim is to compare two interior vertices u and v of \mathcal{T} . Was u more likely before (of lower rank than) v or v before u ? In other words, what is the probability

$$\mathbb{P}_{u < v} := \mathbb{P}[r(u) < r(v)]$$

where $r(\mathcal{T})$ is the set of all possible rank functions on \mathcal{T} . Note that it does not hold $\mathbb{P}[r(u) < r(v)] = \mathbb{P}[r(u) > r(v)]$ even with the uniform distribution on the rank functions. The probability $\mathbb{P}_{u < v}$ is equivalent to counting all the possible rank functions on \mathcal{T} in which u has lower rank than v and divide that number by all possible rank functions on \mathcal{T} . One idea is to sum up the probabilities $\mathbb{P}[r(u)=i, r(v)=j]$ in Equation 3 for all $i < j$ which yields to a runtime of $O(|V|^4)$. The following algorithm COMPARE solves the problem in quadratic time. In the following, for a vertex v , the subtree \mathcal{T}_v of \mathcal{T} consists again of v and all its descendants.

Algorithm: COMPARE (\mathcal{T}, u, v)

Input: A rooted binary phylogenetic tree \mathcal{T} and two distinct interior vertices u and v .

Output: The probability $\mathbb{P}_{u < v} := \mathbb{P}[r(u) < r(v) | \mathcal{T}]$.

- 1: Denote the most recent common ancestor of u and v by ρ_1 .
- 2: **if** $\rho_1 = v$ **then**
- 3: **RETURN** $\mathbb{P}_{u < v} = 0$.

```

4: end if
5: if  $\rho_1 = u$  then
6:   RETURN  $\mathbb{P}_{u < v} = 1$ .
7: end if
8: Let  $\mathcal{T}_{\rho_1}$  be the subtree of  $\mathcal{T}$  which is induced
   by  $\rho_1$ .
9: Delete the vertex  $\rho_1$  from  $\mathcal{T}_{\rho_1}$ . The two
   evolving subtrees are labeled  $\mathcal{T}_u$  and  $\mathcal{T}_v$  with
    $u \in \mathcal{T}_u$  and  $v \in \mathcal{T}_v$ .
10: Run RANKPROB( $\mathcal{T}_u, u$ ) and RANKPROB( $\mathcal{T}_v, v$ )
   to get  $\mathbb{P}[r(u) = i]$  on  $\mathcal{T}_u$  and  $\mathbb{P}[r(v) = i]$  on  $\mathcal{T}_v$ 
   for all possible  $i$ .
11: for  $i = 1, \dots, |\dot{V}_{\mathcal{T}_u}|$  do
12:    $ucum(i) := \sum_{k=1}^i \mathbb{P}[r(u) = k]$ 
13: end for
14:  $\mathbb{P}_{u < v} = 0$ 
15: for  $i = 1, \dots, |\dot{V}_{\mathcal{T}_v}|$  do
16:   for  $j = 1, \dots, |\dot{V}_{\mathcal{T}_u}|$  do
17:      $p := \mathbb{P}[r(v) = i] \cdot \binom{i-1+j}{j}$ 
        $\cdot \binom{|\dot{V}_{\mathcal{T}_v}|-i+|\dot{V}_{\mathcal{T}_u}|-j}{|\dot{V}_{\mathcal{T}_u}|-j} \cdot ucum(j)$ 
18:      $\mathbb{P}_{u < v} := \mathbb{P}_{u < v} + p$ 
19:   end for
20: end for
21:  $tot := \binom{|\dot{V}_{\mathcal{T}_u}|+|\dot{V}_{\mathcal{T}_v}|}{|\dot{V}_{\mathcal{T}_v}|}$ 
22:  $\mathbb{P}_{u < v} := \mathbb{P}_{u < v} / tot$ 
23: RETURN  $\mathbb{P}_{u < v}$ 

```

Theorem 8. *The algorithm COMPARE returns the value*

$$\mathbb{P}_{u < v} = \mathbb{P}[r(u) < r(v)].$$

The runtime of COMPARE is $O(|\dot{V}|^2)$.

Proof. Note that the probability of u having smaller rank than v in tree \mathcal{T}_{ρ_1} equals the probability of u having smaller rank than v in tree \mathcal{T} , since for any rank function on \mathcal{T}_{ρ_1} , there is the same number of linear extensions to get a rank function on the tree \mathcal{T} .

So it is sufficient to calculate the probability $\mathbb{P}_{u < v}$ in \mathcal{T}_{ρ_1} . If $\rho_1 = u$ then u is an ancestor of v in

\mathcal{T} , so return $\mathbb{P}_{u < v} = 1$. If $\rho_1 = v$ then v is an ancestor of u in \mathcal{T} , so return $\mathbb{P}_{u < v} = 0$.

Now assume that $\rho_1 \neq u$ and $\rho_1 \neq v$. The run of RANKPROB calculates the probability $\mathbb{P}[r(u) = i]$ in the tree \mathcal{T}_u and $\mathbb{P}[r(v) = i]$ in \mathcal{T}_v for all i . Next, combine those two linear orders. Assume that $r(v) = i$ and that j vertices of \mathcal{T}_u are inserted before v . Inserting j vertices of \mathcal{T}_u into the linear order of \mathcal{T}_v before v is possible in $\binom{i-1+j}{j}$ ways (see Remark 1). Putting the remaining vertices in a linear order is possible in $\binom{|\dot{V}_{\mathcal{T}_v}|-i+|\dot{V}_{\mathcal{T}_u}|-j}{|\dot{V}_{\mathcal{T}_u}|-j}$ ways. The probability that the vertex u is among the j vertices which have smaller rank than v is $\mathbb{P}[r(u) \leq j] = ucum(j)$. There are $|r(\mathcal{T}_u)|$ possible linear orders on \mathcal{T}_u and $|r(\mathcal{T}_v)|$ possible linear orders on \mathcal{T}_v . The number of linear orders where vertex v has rank i in \mathcal{T}_v , v has rank $i+j$ in \mathcal{T}_{ρ_1} and $r(u) < i+j$ therefore equals

$$p'_{i,j} = \mathbb{P}[r(v) = i] \cdot |r(\mathcal{T}_v)| \cdot \binom{i-1+j}{j} \cdot \binom{|\dot{V}_{\mathcal{T}_v}|-i+|\dot{V}_{\mathcal{T}_u}|-j}{|\dot{V}_{\mathcal{T}_u}|-j} \cdot ucum(j) \cdot |r(\mathcal{T}_u)|$$

Adding up the p' for each i and j gives the number of linear orders where u has smaller rank than v .

Combining a linear order on \mathcal{T}_v with a linear order on \mathcal{T}_u is possible in

$$tot := \binom{|\dot{V}_{\mathcal{T}_u}|+|\dot{V}_{\mathcal{T}_v}|}{|\dot{V}_{\mathcal{T}_v}|}$$

different ways (see Remark 1). There are $|r(\mathcal{T}_u)|$ linear orders on \mathcal{T}_u and $|r(\mathcal{T}_v)|$ linear orders on \mathcal{T}_v , so on \mathcal{T}_{ρ_1} , there are

$$tot' := \binom{|\dot{V}_{\mathcal{T}_u}|+|\dot{V}_{\mathcal{T}_v}|}{|\dot{V}_{\mathcal{T}_v}|} |r(\mathcal{T}_v)| |r(\mathcal{T}_u)|$$

linear orders. Therefore:

$$\mathbb{P}_{u < v} = \frac{\sum_{i,j} p'_{i,j}}{tot'} = \frac{\sum_{i,j} p_{i,j}}{tot}$$

with

$$p_{i,j} = \mathbb{P}[r(v) = i] \cdot \binom{i-1+j}{j} \cdot \binom{|\dot{V}_{\mathcal{T}_v}|-i+|\dot{V}_{\mathcal{T}_u}|-j}{|\dot{V}_{\mathcal{T}_u}|-j} \cdot ucum(j).$$

This shows that COMPARE works correct.

Since RANKPROB has quadratic runtime, COMPARE also has quadratic runtime.

Acknowledgements

We thank Arne Mooers for very helpful comments and suggestions on earlier versions of this manuscript and the two anonymous referees for a very careful report.

The Second author's work is partially supported by grant NSF-DMS-0241246

References

- Aldous, D.J. 2001. Stochastic models and descriptive statistics for phylogenetic trees, from Yule to today. *Statist. Sci.*, 16(1):23–34. ISSN 0883-4237.
- Baum, B.R. 1992. Combining trees as a way of combining data sets for phylogenetic inference, and the desirability of combining gene trees. *Taxon*, 41(1):3–10.
- Brown, J.K.M. 1994. Probabilities of evolutionary trees. *Syst. Biol.*, 43(1):78–91.
- Edwards, A.W.F. 1970. Estimation of the branch points of a branching diffusion process. (With discussion.). *J. Roy. Statist. Soc. Ser. B.*, 32:155–174. ISSN 0035-9246.
- Gernhard, T. 2006. Stochastic models of speciation events in phylogenetic trees. Diplom thesis.
- Harding, E.F. 1971. The probabilities of rooted tree-shapes generated by random bifurcation. *Advances in Appl. Probability*, 3:44–77. ISSN 0001-8678.
- Hey, J. 1992. Using phylogenetic trees to study speciation and extinction. *Evolution*, 46:627–640.
- Nee, S.C. and May, R.M. 1997. Extinction and the loss of evolutionary history. *Science*, 278:692–694.
- Nordborg, M. 2001. Coalescent theory. *Handbook of Statistical Genetics*, p179–212.
- Page, B. 1991. Random cladograms and null hypotheses in cladistic biogeography. *Systematic Zoology*, 40:54–62.
- Pinelis, I. 2003. Evolutionary models of phylogenetic trees. *Roy. Soc. Lond. Proc. Ser. Biol. Sci.*, 270(1522):1425–1431+15. ISSN 0962-8452. With an electronic appendix [DOI 10. 1098 spb. 2003. 2374].
- Ragan, M. 1992. Phylogenetic inference based on matrix representation of trees. *Mol. Phylogenet. Evol.*, 1:53–58.
- Semple, C. and Steel, M. 2003. *Phylogenetics*, volume 24 of *Oxford Lecture Series in Mathematics and its Applications*. Oxford University Press, Oxford. ISBN 0-19-850942-1.
- Vos, R.A. and Mooers, A.O. A new dated supertree of the primates. *Systematic Biology, in Revision*.
- Yule, G.U. 1924. A mathematical theory of evolution: based on the conclusions of Dr. J.C. Willis. *Philos. Trans. Roy. Soc. London Ser. B.*, 213:21–87.
- Zhaxybayeva, O.D. and Gogarten, J.P. 2004. Cladogenesis, coalescence and the evolution of the three domains of life. *Trends in Genetics*, 20:182–187.