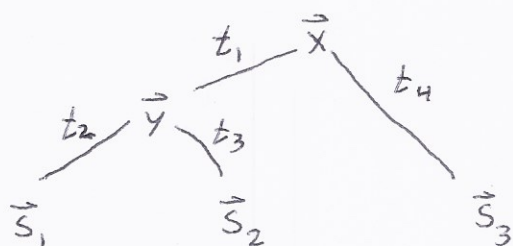


Let's imagine we have some gene sequences  $\vec{S}_1$ ,  $\vec{S}_2$ , and  $\vec{S}_3$ . If they evolved without recombination, then there is some true evolutionary relationship among them that can be represented by a phylogenetic tree!



Here is one possible tree topology. The number of possible topologies increases rapidly with the number of sequences

The tree,  $T$ , is fully specified by the topology and the branch lengths ( $t_1, t_2, t_3, t_4$ ).

Let's say we have some evolutionary model that tells us ~~the probability that~~ ~~each~~ each sequence  $\vec{S}_1$  is likely to evolve to some other sequence  $\vec{S}_2$  after some duration of time  $t$ :

$$Pr(\vec{S}_2 | \vec{S}_1, t) \leftarrow \begin{array}{l} \text{probability of being } \vec{S}_2 \text{ after duration} \\ t \text{ given that sequence was } \vec{S}_1 \\ \text{at start of time duration} \end{array}$$

Note that this transition probability does not depend on the absolute date, only the duration of time  $t$ . So we are assuming evolution is time homogenous. When might this be untrue?

Let's return to the tree above. We would like to write the likelihood of the sequence data given the tree:

$$Pr(\vec{S}_1, \vec{S}_2, \vec{S}_3 | T) = \sum_{\vec{x}} \sum_{\vec{y}} Pr(\vec{x}) \cdot Pr(\vec{S}_3 | \vec{x}, t_4) \cdot Pr(\vec{y} | \vec{x}, t_1) \cdot Pr(\vec{S}_1 | \vec{y}, t_2) \cdot Pr(\vec{S}_2 | \vec{y}, t_3)$$

Felsenstein pruning algorithm, a form of dynamic programming  $\rightarrow$   $= \sum_{\vec{x}} Pr(\vec{x}) \cdot Pr(\vec{S}_3 | \vec{x}, t_4) \cdot \sum_{\vec{y}} Pr(\vec{y} | \vec{x}, t_1) \cdot Pr(\vec{S}_1 | \vec{y}, t_2) \cdot Pr(\vec{S}_2 | \vec{y}, t_3)$

The  $Pr(\vec{S} | \vec{x}, t)$  terms are the transition probabilities given by our evolutionary model.

What is  $Pr(\vec{x})$ ? It is the root frequencies or stationary state.

This is the principal eigenvector of the transition matrix.

We now have an equation to compute:

$$\Pr(\text{sequences} | \text{tree} = T, \text{model} = \Pr(\vec{s}_1 | \vec{s}_2, t))$$

Maximum likelihood:

Find tree that maximizes

$$\Pr(\text{sequences} | \text{tree}, \text{model})$$

Bayesian:

Place prior over trees:  $\Pr(\text{tree})$

Then compute posterior over trees:

$$\Pr(\text{tree} | \text{sequences}) \propto \Pr(\text{tree}) \cdot \Pr(\text{sequences} | \text{tree})$$

Computation uses MCMC

Choice of method depends on taste, computational difficulty, and question at hand.

Sometimes a researcher just wants a tree to use as figure.

Other times they want to compute posterior distribution over a statistic calculated from tree.