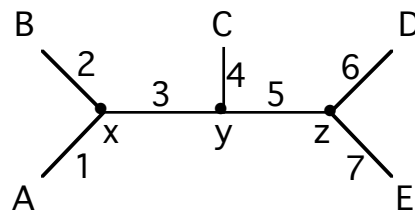


Lecture 4 – Introduction to Trees and Optimality Criteria: Parsimony

I. To begin our treatment of estimating phylogeny, I want to first introduce some generalities about trees, so that we'll have a common starting point, and then discuss the various optimality criteria that one may choose for evaluating trees.

A. Some Basic Terms. There are two sets of terms that are commonly used: one by mathematicians and one by systematists.

In a sense, the mathematicians' terms are more appropriate because phylogenetic trees fall into graph theory, however, having been trained as a systematist, I think in those terms.



1. Branches: If there are n -taxa in an unrooted tree, there are $2n-3$ branches, labeled 1 – 7.

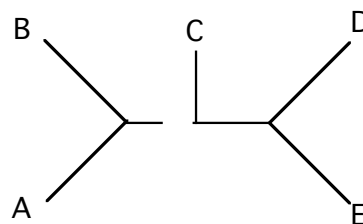
- 1, 2, 4, 6, & 7 are external, or terminal branches (leaves in math jargon)
- 3 & 5 are internal branches (edges)

2. Nodes are points where branches intersect or terminate

Nodes A – E are terminals, or terminal nodes, whereas nodes x, y, & z are internal nodes.

3. We can denote trees in parenthetical notation (sometimes called New Hampshire or Newick format). This would be represented by $((A,B),C,(D,E))$.

4. We can break trees into sub-trees. If we break branch 3, we have two sub-trees (A,B) and (C,(D,E)).



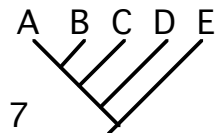
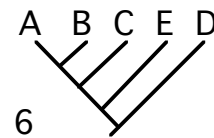
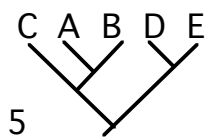
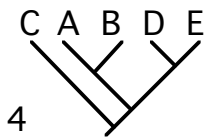
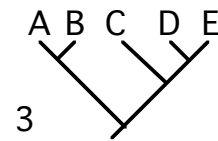
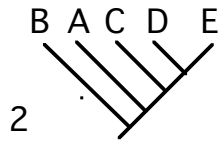
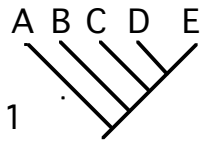
B. Rooting – The tree above is an unrooted tree.

This implies a lack of directionality, such that one may traverse the tree in any direction one finds convenient and start the traversal at any point.

Modern phylogenetic analyses usually estimate unrooted trees, and the root is usually determined by the inclusion of an outgroup. The point at which the outgroup connects with the ingroup is the root.

This is usually included in the estimation, not simply added after an ingroup topology is estimated (as I'll do in this demonstration).

All of the rooted trees below can be derived from the unrooted tree by placing the root (common ancestor) on each branch. These rooted trees *do* imply directionality; nodes deeper in the tree represent divergence events that happened earlier. This illustrates how rapidly the number of alternative trees grows as one adds taxa. Addition of a single node (the root node) resulted in the seven trees below being derivable from the single tree above.



C. The scope of the phylogeny problem.

One other important point to make right off can be seen by the rooting exercise.

The number of possible trees increases incredibly fast as the number of taxa increases.

We can think of rooting as adding another taxon to this tree. There are 7 places we can add the 6th taxon and still be congruent with this 5-taxon tree, and there are 15 possible 5-taxon trees.

What this means is that the number of possible alternative trees will quickly be far too vast for us to examine every possible tree for a phylogenetic question involving very many taxa.

Number of possible trees for N taxa:

$$B(N) = \prod_{i=3}^N (2i - 5)$$

Number of Taxa	Number of Unrooted Trees
3	1
4	3
5	15
6	105
7	945
8	10,395
9	135,135
10	2.027 X 10 ⁶
22	3 X 10 ²³
50	3 X 10 ⁷⁴
100	2 X 10 ⁸²
1000	2 X 10 ^{2,860}
10,000,000	5 X 10 ^{68,667,340}

So the task of phylogeny inference is quite formidable, which is why the idea that phylogenetic trees can be “made easy” makes little sense (Hall, B. 2001. *Phylogenetic Trees Made Easy*, Sinauer, Sunderland, MA).

Some approaches to phylogeny estimation don't try at all to search among these alternatives, but rely on a series of rules to build one. These are referred to as algorithmic methods.

Other approaches actually try to search this vast tree space in an effort to find the best tree under some optimality criterion.

II. Parsimony – Before we deal with how to choose from among all these trees, we first need to know how to evaluate any one tree.

Both algorithmic and optimality approaches rely on being able to evaluate a particular tree.

This is the most frequently used criterion. It's certainly the easiest to compute, it's also (remember) the only philosophically acceptable criterion for some systematists.

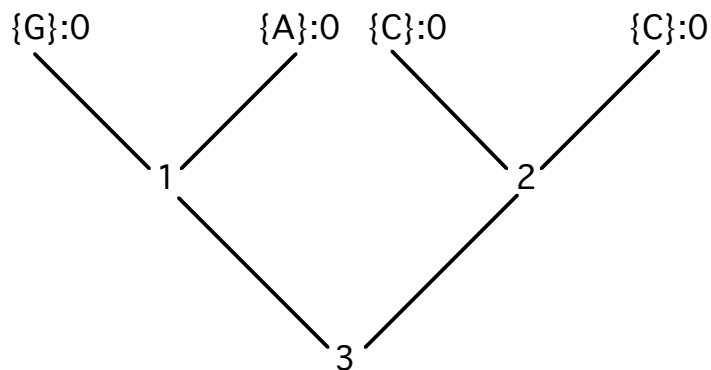
First, the parsimony score of a tree (i.e., its length, L_τ) for the entire data set is given by:

$$L_\tau = \sum_{i=1}^s w_i l_i$$

That is, it's the weighted sum of the character lengths, l_i , across all s characters. For now, we'll consider equal weights (i.e., all $w_i = 1$)

There are two ways one may calculate the L_τ 's

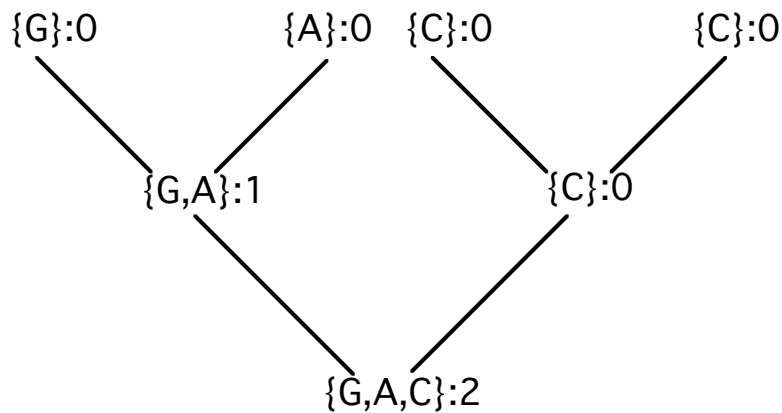
A. The Fitch Algorithm applies to matrix with unordered characters and all character-state transformations equal. It operates on **state sets** and **accumulated lengths**.



There are four terminals for this (arbitrarily rooted) tree: one with G, one with A and two with C. We erect a state set at each terminal node and assign an accumulated length of zero to terminal nodes. This is the minimum number of changes in the daughter subtree.

We then proceed down the tree to internal nodes and calculate the state sets for each, as follows:

- 1 – Form the intersection of the two daughter nodes. If the intersection is non-empty, assign the set for the internal node equal to the intersection. The accumulated length of the internal node is the sum of those of the daughter nodes.
- 2 - If the intersection is empty, we assign the union of the two daughter nodes to the state set for the internal node. The accumulated length is the sum of those of the daughter nodes plus one.



Therefore, this character has a minimum length of 2 steps for this tree. In addition to determining the length of this tree for this character, we've also assigned possible ancestral states at each of the nodes.

Remember that this is only applicable to characters for which transformations among all characters are unordered and have the same cost.

2. Sankoff's Algorithm is more general, and is applicable for any type of character.

Let's assume that we think the following transformation matrix is appropriate:

	A	C	G	T
A	--	4	1	4
C	4	--	4	1
G	1	4	--	4
T	4	1	4	--

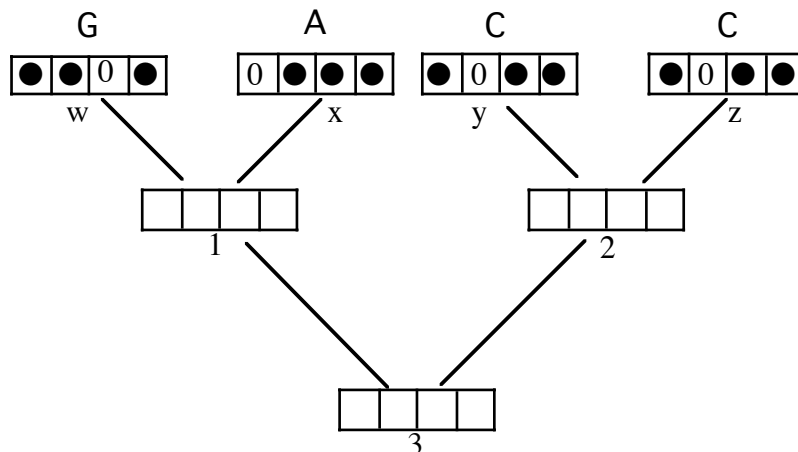
The convention is to use an alphabetical order, and c_{ij} is the cost (in # of steps) to transform from state i to state j . It is common practice in molecular phylogeny to use a non-uniform matrix (called a step matrix), and this one is symmetric. This step matrix implies that transversions are expected to be ca. 4x as reliable as transitions. More (much more) about this later.

Following the Sankoff algorithm, we erect a **character-state vector** at each node, and find the cost of assigning each possible state at each node.

For each node k , we calculate the cost associated with each possible character-state assignment to that node, $s_k(i)$ this is analogous to the assumed cost above, in that it represents the cost of all changes leading to the daughter nodes, given assignment of state i to node k .

Again, we start at the tips, and the terminal vectors are filled with either 0, if the state is present or a dot (or infinity, or an asterisk...), if the state is absent in that taxon.

Since we know states at the terminal nodes, we only need worry about the cost of transformations to the state present. Therefore, the vector at node 1 is filled in as such:



$$s_{1(A)} = c_{AG} + c_{AA} = 1 + 0 = 1,$$

$$s_{1(C)} = c_{CG} + c_{CA} = 4 + 4 = 8,$$

$$s_{1(G)} = c_{GG} + c_{GA} = 0 + 1 = 1$$

$$s_{1(T)} = c_{TG} + c_{TA} = 4 + 4 = 8$$

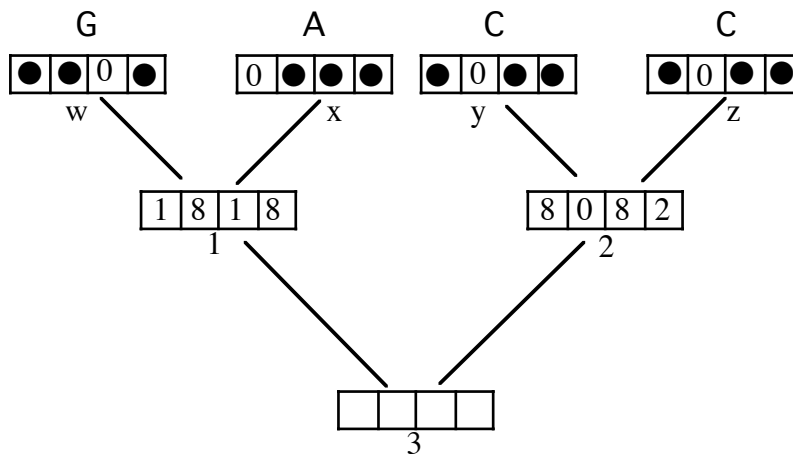
Similarly, the node 2 character state vector is as follows:

$$s_{2(A)} = 4 + 4 = 8$$

$$s_{2(C)} = 0 + 0 = 0$$

$$s_{2(G)} = 4 + 4 = 8$$

$$s_{2(T)} = 1 + 1 = 2$$



Now, since we don't have sequences for nodes 1 & 2, when we calculate the cells of the node 3 vector, we have to consider all possible character states for each.

Thus, the cells of the character-state vector incorporate the minimum values for each possible reconstruction of its immediate daughter nodes.

$$s_{3(A)} = \min[s_{1Aj} + c_{Aj}] + \min[s_{2Aj} + c_{Aj}]$$

$$= \min [1,12,2,12] + \min[8,4,9,6] = 1 + 4 = 5$$

$$s_{3(C)} = \min[s_{1Cj} + c_{Cj}] + \min[s_{2Cj} + c_{Cj}]$$

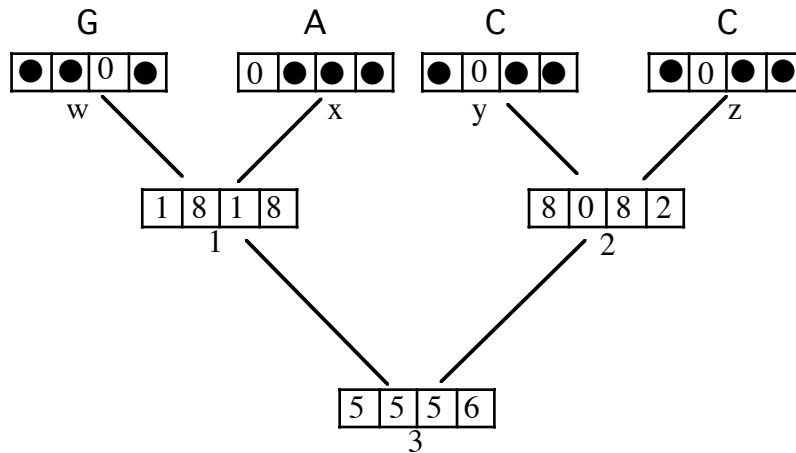
$$= \min [5,8,5,9] + \min[12,0,12,3] = 5 + 0 = 5$$

$$s_{3(G)} = \min[s_{1Gj} + c_{Gj}] + \min[s_{2Gj} + c_{Gj}]$$

$$= \min [2,12,1,12] + \min[9,4,8,6] = 1 + 4 = 5$$

$$s_{3(T)} = \min[s_{1Tj} + c_{Tj}] + \min[s_{2Tj} + c_{Tj}]$$

$$= \min [5,9,5,8] + \min[12,1,12,2] = 5 + 1 = 6$$



Now, the node-3 character-state vector is defined and all the nodes have been addressed:

Thus, this character has a minimum cost of five steps under a 4:1 weighting of transversions.

We can do this for all characters on this tree and a weighted sum taken across characters to compute the parsimony score of this tree.

Points to note:

- 1) Two types of weighting are possible: weighting of transformations within characters (which we demonstrated with the step matrix) and weighting among characters, which are reflected in the weighted sum of lengths across characters.
- 2) One can't compare tree lengths across weighting schemes. In the first example with all transformations having the same cost, the length of the character on this tree was 2. In the second, with a 4:1 step matrix to weight transversions more heavily, the length was 5. This prevents us from objectively comparing weighting schemes.