Introduction to Bioinformatics

Phylogenetic trees

# Inferring the Past: Phylogenetic Trees

- p The biological problem
- p Parsimony and distance methods
- p Models for mutations and estimation of distances
- p Maximum likelihood methods

### Phylogeny

- We want to study ancestordescendant relationships, or *phylogeny*, among groups of organisms
- P Groups are called *taxa* (singular: *taxon*)
- P Organisms are usually called *operational taxonomic units* or *OTUs* in the context of phylogeny



### Phylogenetic trees

- Leaves (external nodes)
  ~ species, observed (OTUs)
- p Internal nodes ~ ancestral species/divergence events, not observed
- p Unrooted tree does not specify ancestordescendant relationships beyond the observation "leaves are not ancestors"



Unrooted tree with 5 leaves and 3 internal nodes.

*Is node 7 ancestor of node 6?* 

### Phylogenetic trees

 $R_2$ 

 $R_2$ 

root(R2)

**R**₁

6′

tootlen

- Rooting a tree specifies all ancestor-descendant relationships in the tree
- P Root is the ancestor to the other species
- P There are n-1 ways to root
  a tree with n nodes
  R<sub>1</sub>

#### Questions

- p Can we enumerate all possible phylogenetic trees for *n* species (or sequences?)
- p How to score a phylogenetic tree with respect to data?
- p How to find the best phylogenetic tree given data?

# Finding the best phylogenetic tree: naive method

- p How can we find the phylogenetic tree that best represents the data?
- p Naive method: enumerate all possible trees
- p How many different trees are there of n species?
- p Denote this number by b<sub>n</sub>

#### Enumerating unordered trees

p Start with the only unordered tree with 3 leaves  $(b_3 = 1)$ 



p Consider all ways to add a leaf node to this tree



- Fourth node can be added to 3 different branches (edges), creating 1 new internal branch
- p Total number of branches
  is n external and n 3
  internal branches
- P Unrooted tree with n
  leaves has 2n 3 branches

#### Enumerating unordered trees

p Thus, we get the number of unrooted trees

$$o_n = (2(n - 1) - 3)b_{n-1} = (2n - 5)b_{n-1}$$
  
=  $(2n - 5) * (2n - 7) * ... * 3 * 1$   
=  $(2n - 5)! / ((n-3)!2^{n-3}), n > 2$ 

p Number of rooted trees b'n is

that is, the number of unrooted trees times the number of branches in the trees

# Number of possible rooted and unrooted trees

n	B <sub>n</sub>	b'n
3	1	3
4	3	15
5	15	105
6	105	945
7	954	10395
8	10395	135135
9	135135	2027025
10	2027025	34459425
20	2.22E+020	8.20E+021
30	8.69E+036	4.95E+038

#### Too many trees?

p We can't construct and evaluate every phylogenetic tree even for a smallish number of species

#### p Better alternative is to

- n Devise a way to evaluate an individual tree against the data
- n Guide the search using the evaluation criteria to reduce the search space