

## Distance Methods

	A	B	C	D
A	0	10	6	4
B	10	0	12	8
C	6	12	0	6
D	4	8	6	0

1. Apply the UPGMA algorithm to the matrix to find the UPGMA tree.
2. Assume that the first step of the neighbor-joining algorithm joins taxa A and C, that the edge to A has length 0.5, the edge to C has length 5.5, and that the new distance matrix relating the remaining groups is as follows.

	A/C	B	D
A/C	0	8	2
B	8	0	8
D	2	8	0

Find the neighbor-joining tree.

3. For both the UPGMA tree and the neighbor-joining tree, find the pairwise distance matrix that agrees with the tree perfectly. Compare these with the original distance matrix. Which of the two estimates a tree where the distances are closest to the original?

## Maximum Likelihood

4. Answer true or false: *There is an algorithm similar to that for parsimony which allows rapid calculation of the log-likelihood of a tree topology.*
5. Answer true or false: *For a specified likelihood model, there is an algorithm which allows rapid calculation of the log-likelihood of a tree (where both the topology and branch lengths are specified).*
6. Answer true or false: *All likelihood models, such as Jukes-Cantor and the General Time Reversible model lead to the same maximum likelihood tree topology, but the branch lengths may differ.*
7. Answer true or false: *For a specified maximum likelihood model, the maximum likelihood tree is the one that makes the probability of the observed data as high as possible.*
8. Recall that the AIC score is  $-2(\log\text{-likelihood}) + 2(\# \text{ of parameters})$  and that the AIC criterion is to choose the model that minimizes AIC. Which model would be preferred from this list according to AIC?

Model	$\ln L$	# of parameters
1	-3068	25
2	-2953	25
3	-2953	26
4	-2935	28
5	-2680	28
6	-2616	29

## Bootstrap

9. Suppose a molecular alignment includes 300 sites for seven taxa. The maximum-likelihood tree from this data contains a clade with taxa 1–3. Explain how the bootstrap could be used to generate 1000 bootstrap data sets to quantify the support for this clade.
10. What is the probability that the first site is not included in the first bootstrap data set?
11. For data sets with  $n$  columns where  $n = 100, 200, 300, \dots, 1000$ , find the probability that the first site is not part of the first bootstrap data set. What happens to these numbers as  $n$  increases?
12. For a large data set, estimate to the nearest 10% the proportion of bootstrap data sets you would expect to include any single specific site.
13. Assume a molecular data set for 4 taxa and the maximum parsimony method of tree estimation. Suppose that 5 sites favor the tree with an A/B clade, 3 sites favor the tree with an A/C clade, 2 sites favor a tree with an A/D clade, and the remaining 90 sites are not parsimony informative.
  - (a) What is the probability that there are no parsimony informative sites in a bootstrap data set?
  - (b) For each bootstrap data set, how would you determine the parsimony tree?
  - (c) What potential difficulty would need to be addressed to find the bootstrap support for the most parsimonious tree?

## Bayesian Inference

14. Answer true or false: *Similar to maximum parsimony and maximum likelihood, Bayesian inference requires a search through tree-space for the best tree.*
15. Answer true or false: *Similar to maximum parsimony and maximum likelihood, the bootstrap is often used to assess uncertainty in the Bayesian tree.*
16. Answer true or false: *The Bayesian approach includes both likelihood models of molecular evolution and prior probability distributions for the tree and model parameters.*
17. In an inference problem with four taxa and three possible unrooted tree topologies, assume that it is possible to calculate the likelihood of data for each tree topology by averaging over some prior distribution on branch lengths and parameters of the likelihood model. The following table shows these probabilities on the natural log scale.

Tree	$\ln L$
1	-5000
2	-5002
3	-5010

In other words, the probability of the observed data (averaging over uncertainty in other parameters) is  $\exp(-5000)$  for the first tree topology and so on.

- (a) If a person places a uniform prior probability distribution on the three tree topologies (and uses the other unspecified prior distributions on branch lengths and other model parameters), what will the Bayesian posterior distribution for these three trees be? (*Hint: the problem is similar to the box and balls problem from lecture — drawing a tree diagram can help.*)
- (b) If the prior distribution on tree topologies has instead been (0.02, 0.49, 0.49), which tree would have had the highest posterior probability (and what would it have been)?