

Introduction to Phylogenetics

Week 5

Distance methods

	Charater- based method	Non-character based
Explicit evolutionary model	Maximum likelihood	Pairwise distance
No explicit evolutionary model	Maximum parsimony	

I. Distance Methods

- Try to fit tree to genetic distance (d -distance) matrix
- d -distance determined based on observed distance (p -distance)
- Use with evolution/substitution models – attempt to determine how many sites have actually changed
- Determine topology based on generated d -distance matrix

I. Distance Methods

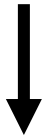
Align sequences



Calculate observed distance (p -distance)



Estimate evolutionary distance (d -distance)



Calculate for every sequence vs sequence

I. Distance Methods

- Requires on accurate calculation d -distance (by default good alignment)
- Dependent on correct choice of evolutionary model
- Trees build using:
 - cluster analysis
 - minimum evolution
- Clustering sensitive to unequal evolutionary rates in different species

II. Cluster Analysis

- Based on taxonomic phenograms
- Applied to create ultrametric trees
 - Uses math based on triangle inequities
 - Ultrametric strengthened triangle inequity
 - Replace 'sides' with 'taxa'

$$d_{AC} \in (d_{AB}, d_{BC})$$

- Two of three distances are equal or larger than a third
- Assumes molecular clock

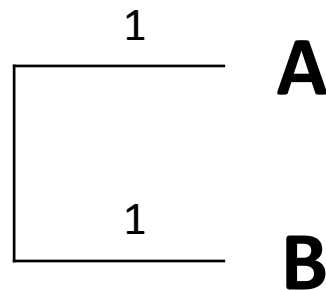
II. Cluster Analysis

- Two primary clustering methods:
 - Unweighted-pair group method with arithmetic mean (UPGMA)
 - Weighted-pair group method with arithmetic mean (WPGMA)
- Trees are built step-wise
- Grouping OTUs base on smallest genetic distance to create new d value
- Calculate distance from new node u to any other node k

WPGMA

$$d_{uk} = \frac{d_{(A,B)k} + d_{Ck}}{2}$$

	A	B	C	D	E
B	2				
C	4	4			
D	6	6	6		
E	6	6	6	4	
F	8	8	8	8	8



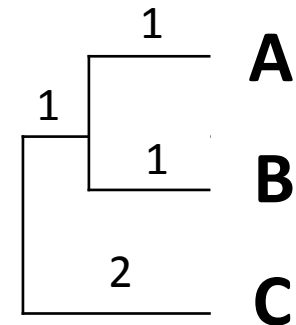
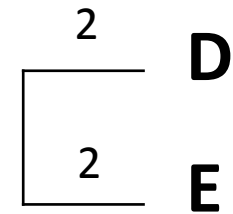
	AB	C	D	E
C	4			
D	6	6		
E	6	6	4	
F	8	8	8	8

$$d_{(AB)C} = \frac{d_{(AC)} + d_{BC}}{2} = 4$$

$$d_{(AB)D} = \frac{d_{(AD)} + d_{BD}}{2} = 6$$

$$d_{(AB)E} = \frac{d_{(AE)} + d_{BE}}{2} = 6$$

$$d_{(AB)F} = \frac{d_{(AF)} + d_{BF}}{2} = 8$$



II. Cluster Analysis (UPGMA)

- Distances averaged on the number of OTUs
- Changes value of u and k

$$d_{uk} = \frac{(N_{AB}d_{(A,B)k} + N_Cd_{Ck})}{(N_{AB} + N_C)}$$

N_{AB} = # OTUs in cluster AB

N_C = # OTUs in cluster C

- Mostly affects data that is not ultrametric

III. Minimum Evolution (ME)

- Clustering methods sensitive to variation
- NJ creates additive distance trees
- Uses four-point metric condition

$$d_{AB} + d_{CD} \leq \max(d_{AC} + d_{BD}, d_{AD} + d_{BC})$$

- Are used to create an unrooted tree – sum of distance between pair OTUs is sum (rather than average) of length connecting them
- Tree algorithm cannot always satisfy four-point condition – tree distances different

III. Minimum Evolution (ME)

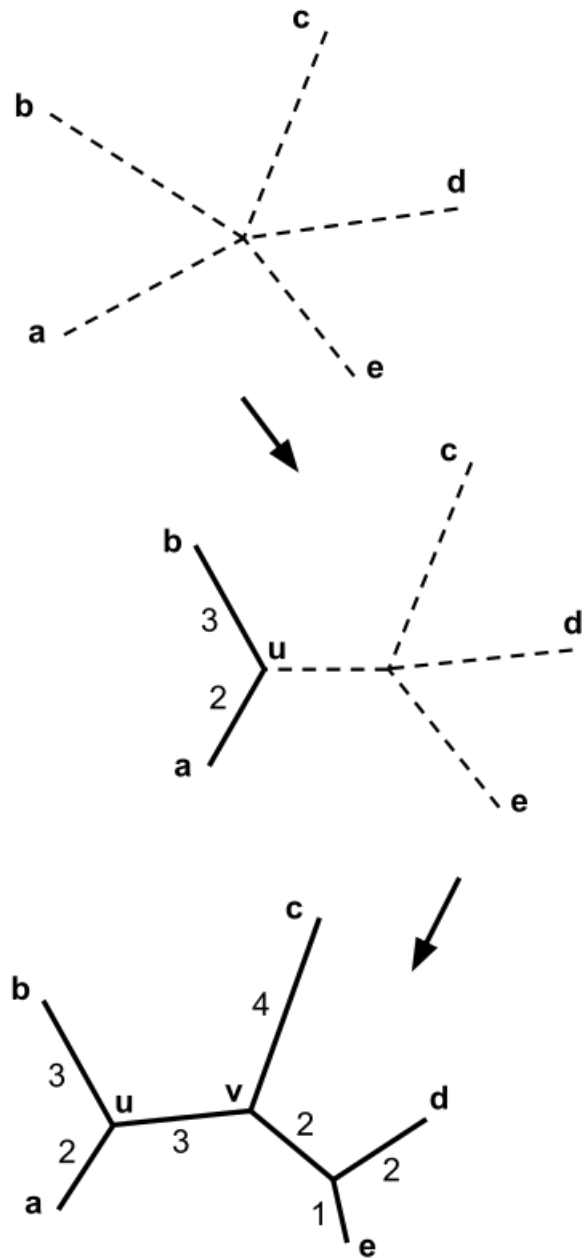
- Clustering methods sensitive to variation
- NJ creates additive distance trees
- Uses four-point metric condition

$$d_{AB} + d_{CD} \leq \min(d_{AC} + d_{BD}, d_{AD} + d_{BC})$$

- Are used to create an unrooted tree – sum of distance between pair OTUs is sum (rather than average) of length connecting them
- Tree algorithm cannot always satisfy four-point condition – tree distances different

	a	b	c	d
b	5			
c	9	10		
d	9	10	8	
e	8	9	7	3

Branches = $2n-3$



III. Minimum Evolution

- ME much better algorithm, but cannot take into account stochastic variation
- Minimum evolution method attempts to minimize branch lengths of tree
- Gives a best estimate of phylogeny

$$S = \sum_{i=1}^{2n-3} v_i$$

n = number of taxa in tree

v_i = i th branch

III. ME – Neighbor-Joining

- Can estimate branch length from distance matrix
- Optimizes so much that might not find best tree (which might be bigger)
- Have to explore many different topologies to find smallest tree - only works for smaller trees
- Can overcome some issues with NJ approach
- NJ does not assume an evolutionary clock

III. ME – Neighbor-Joining

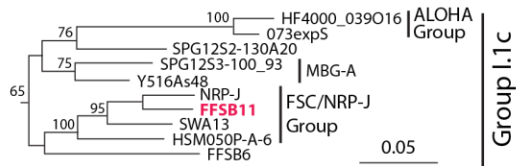
- Most common method to estimate distance trees
- Minimizes S value, but does for each branch added (S not globally minimized)
- Generally generates similar tree to ME method
- Can evaluate final trees using bootstrapping
- Can combine bootstrapping with maximum-likelihood (NJML)

IV. Evaluation of trees: Bootstrapping

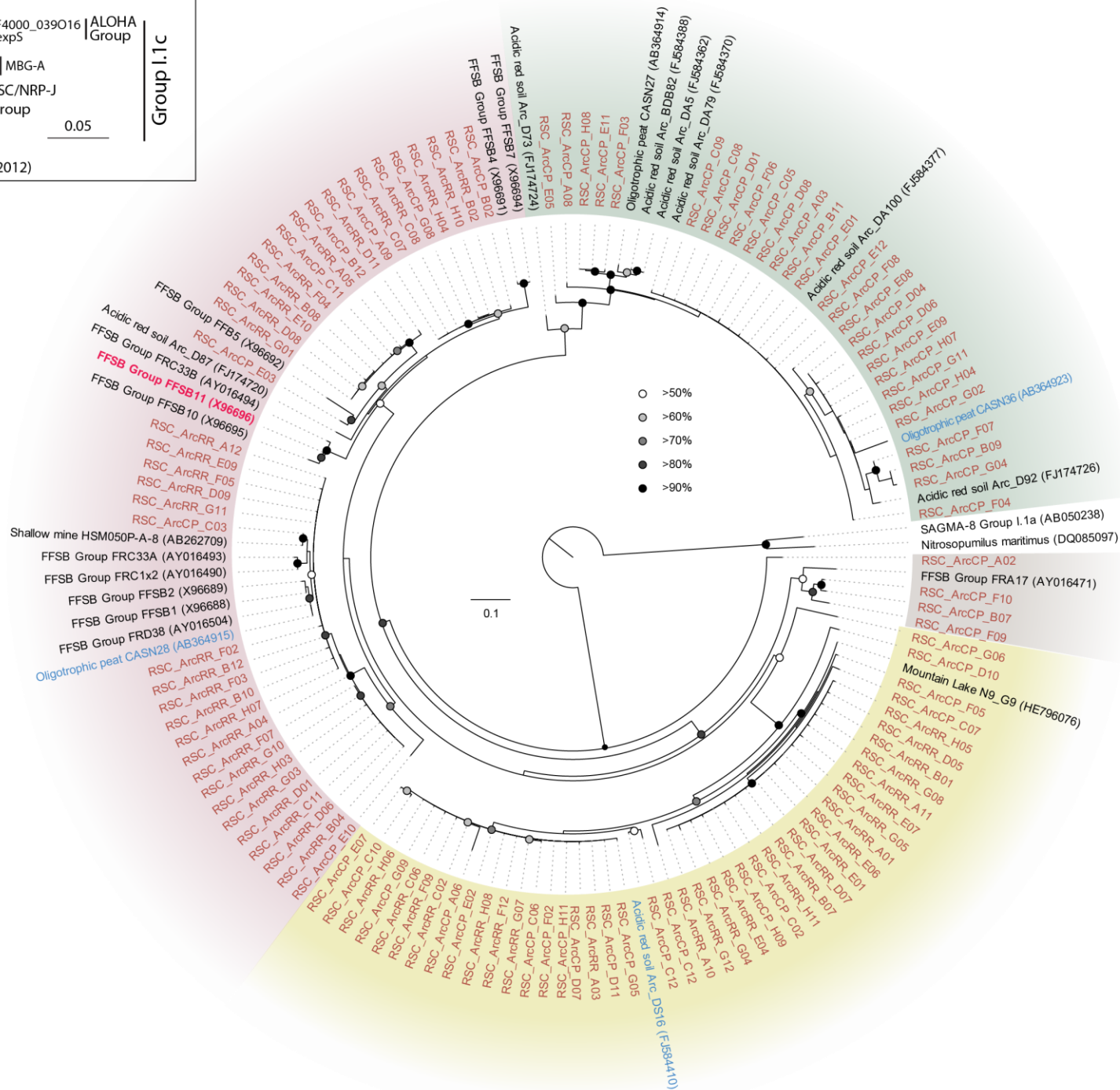
- Calculates statistical error when using an unknown sampling distribution
- Allows approximation underlying distribution
- Resamples original **sequence data** with replacement
- Samples the alignment – creates a tree
- Proportion of clades from bootstrap tree agree with calculated tree
- Provides confidence of grouping/clades calculated

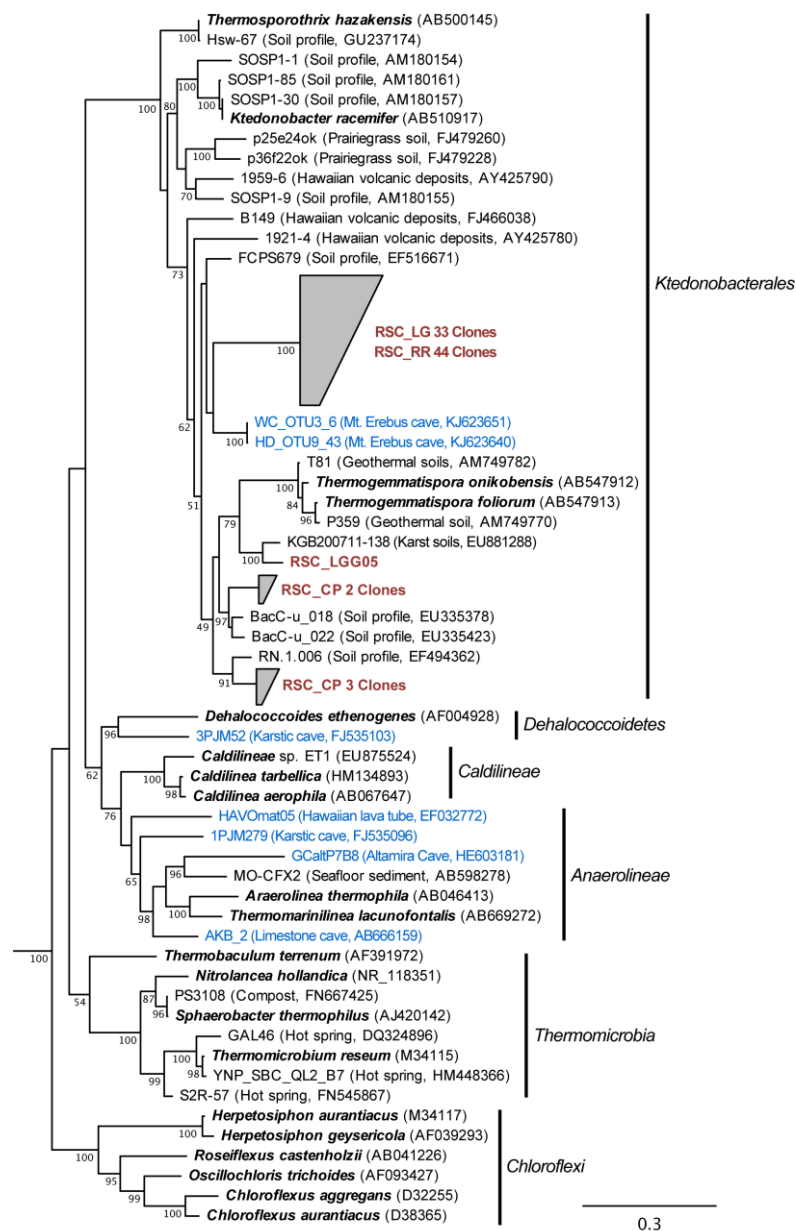
IV. Evaluation of trees: Bootstrapping

- Can be shown as a majority rule consensus
- Can be shown as values on tree
- Bootstrapping demonstrates statistical noise in data
- If your tree is bad, your bootstrap data will not support the topology
 - Poor alignment
 - Wrong substitution model
 - Incorrect assumption distribution



Inset adapted from Durbin and Teske (2012)





IV. Evaluation of trees: Bootstrapping

- Bootstrapping can be fooled by sequences
 - Artificial grouping of data – sequences too similar compared with other sequence data (high G+C)
- Long-branch attraction
 - Incorrect sequence data or rapidly evolving sequences – distant sequences will group together
 - Divergent sequences drawn toward the root

IV. Evaluation of trees: Jackknifing

- Also a subsampling technique
- Deletes random sequences from alignment – recreates tree
- Looks at number of trees match best-tree topology
- Gives a jackknife value that represents % of trees that match clade topology

V. PAUP* and PHYLIP

- Traditionally used for distance/NJ methods
- PAUP* is command-line only – costs \$\$
- PHYLIP free, but very difficult to use due to disjointed program development

<http://evolution.genetics.washington.edu/phylip.html>

- PHYLIP gateway now on Pasteur Server
- <http://mobyte.pasteur.fr/cgi-bin/portal.py#welcome>

V. PHYLIP – Programs

DISTMAT Calculates evolutionary distance (d) values between sequences. Generates distance matrix for tree-building distance methods.

DNADIST Computes four different distances between species from nucleic acid sequences. The distances can then be used in the distance matrix programs. The distances are the Jukes-Cantor formula, one based on Kimura's 2- parameter method, the F84 model used in DNAML, and the LogDet distance.

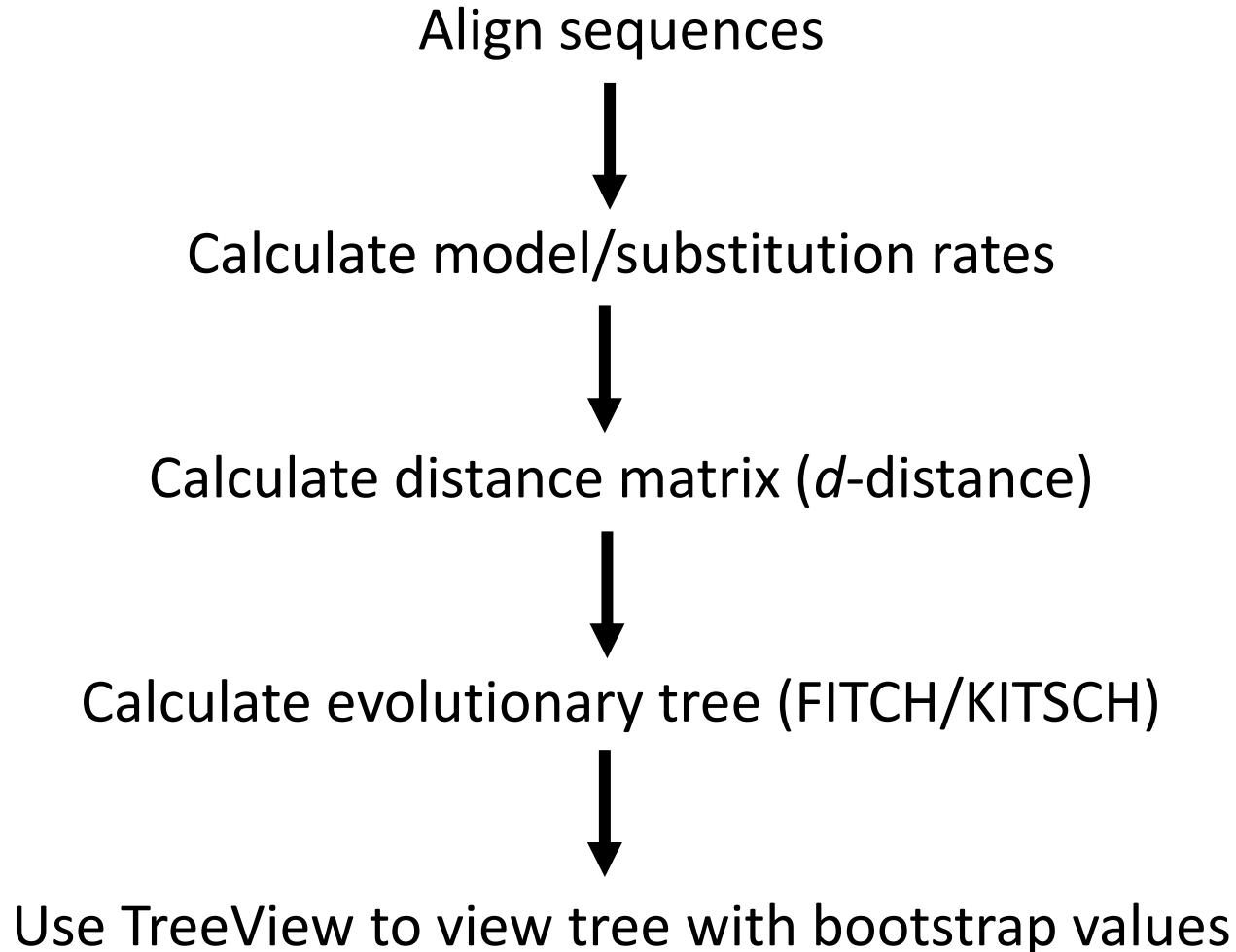
PROTDIST Computes a distance measure for protein sequences, using maximum likelihood estimates based on the Dayhoff PAM matrix, the JTT matrix model, the PBM model, Kimura's 1983 approximation to these, or a model based on the genetic code plus a constraint on changing to a different category of amino acid

FITCH Estimates phylogenies from distance matrix data under the additive tree model.

KITSCH Estimates phylogenies from distance matrix data under the ultrametric model.

NEIGHBOR An implementation of the UPGMA (Average Linkage clustering) method. Methods are very fast and thus can handle much larger data sets.

V. PHYLIP Distance Methods



V. PHYLIP Distance Methods

- PHYLIP will provide trees – image quality not great
- If you generate NEWICK format trees, can use FigTree

<http://tree.bio.ed.ac.uk/software/figtree/>

- NEWICK standard tree format in nested parentheses

$$((A,B),(C,D)) \equiv \begin{array}{c} A & & C \\ & \diagdown \quad \diagup & \\ & \text{---} & \\ & \diagup \quad \diagdown & \\ B & & D \end{array}$$

