

distance-based phylogeny construction algorithms: UPGMA & Neighbor-Joining

UPGMA ~ "unweighted pair group method using arithmetic averages"

o OUV = operational taxonomic unit

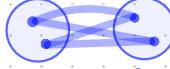
o major assumption for UPGMA: universal biological clock \rightarrow same mutational/evolutionary rate for all species

o input: set of species: $\{A_1, A_2, \dots, A_n\}$

• pairwise distance matrix (D_{ij}) containing distances b/wn species (i.e. alignment scores, Jukes-Cantor dist, etc)

→ Consider 2 clusters of points (or species) S_i and S_j . We want compute d_{ij} : the distance between the 2 clusters

$$d_{ij} = \frac{1}{|S_i| \cdot |S_j|} \sum_{\substack{\text{seqs.} \\ i \in S_i \\ j \in S_j}} d_{seq} = \text{average distance between all pairs of sequences from each cluster}$$

ex:  $\rightarrow d_{ij} = \text{avg of the 4 distances}$

so now does the algorithm use this distance calculation d_{ij} ?

→ suppose S_k is the union of S_i and S_j . S_m is another cluster. what's the distance b/wn S_k and S_m ?

$$d_{km} = \frac{dim|S_i| + dim|S_j|}{|S_i| + |S_j|} \rightarrow \text{now did we get this from our distance formula?}$$

$$dim = \frac{1}{|S_i \cup S_j| \cdot |S_m|} \sum_{\substack{\text{seqs.} \\ i \in S_i \\ j \in S_j \\ m \in S_m}} d_{seq} = \frac{1}{|S_i + S_j| \cdot |S_m|} \left(\sum_{\substack{\text{seqs.} \\ i \in S_i \\ m \in S_m}} d_{seq} + \sum_{\substack{\text{seqs.} \\ j \in S_j \\ m \in S_m}} d_{seq} \right) = \frac{1}{|S_i + S_j| \cdot |S_m|} \cdot |S_i| \sum_{\substack{\text{seqs.} \\ i \in S_i \\ m \in S_m}} d_{seq} + \frac{1}{|S_i + S_j| \cdot |S_m|} \cdot |S_j| \sum_{\substack{\text{seqs.} \\ j \in S_j \\ m \in S_m}} d_{seq}$$

$$= \frac{|S_i|}{|S_i + S_j|} \cdot \frac{1}{|S_m| \cdot |S_i|} \sum_{\substack{\text{seqs.} \\ i \in S_i \\ m \in S_m}} d_{seq} + \frac{|S_j|}{|S_i + S_j|} \cdot \frac{1}{|S_m| \cdot |S_j|} \sum_{\substack{\text{seqs.} \\ j \in S_j \\ m \in S_m}} d_{seq}$$

$$= \frac{|S_i|}{|S_i + S_j|} \cdot dim + \frac{|S_j|}{|S_i + S_j|} dim = \frac{dim|S_i| + dim|S_j|}{|S_i + S_j|}$$

The algorithm:

input: set of sequences $\{A_1, \dots, A_n\}$ and a pairwise distance matrix

output: phylogenetic tree of A_1, \dots, A_n

initialization:

define a cluster for each sequence: $S_k = \{A_k\}$, $1 \leq k \leq n$

define a leaf of T for each input sequence A

iteration: # keep doing this step until you have 2 clusters remaining

Find 2 clusters S_i and S_j for which d_{ij} is minimal

define a new cluster $S_k = S_i \cup S_j$

compute dim for every other cluster m using

define a new node L with daughter nodes i and j . Place L at height $\frac{d_{ij}}{2}$

add L to T and remove i and j

termination: # when there are only 2 clusters remaining $(S_i \text{ and } S_j)$

place root at height $\frac{d_{ij}}{2}$

now, let's turn our attention to the perfect ppt on the website