

### Corrected UPGMA algorithm

Given: a distance matrix D with n species:

1) Initialize n clusters,  $C_1, \dots, C_n$ , each with a single species in it. Create a leaf node for each of the clusters.

2) Define the distance between two clusters as the average pairwise distance between members of the two clusters:

$$d(C_i, C_j) = \frac{\sum_{a \in C_i} \sum_{b \in C_j} D(a, b)}{|C_i| * |C_j|}$$

3) Repeat:

3.1 Pick the two clusters  $C_i$  and  $C_j$  such that  $d(C_i, C_j)$  is minimized.

3.2 Create a new cluster  $C_k = C_i \cup C_j$

3.3 Create a new node in the tree, make it the parent of nodes  $i$  and  $j$ , at height  $d(C_i, C_j)/2$ .

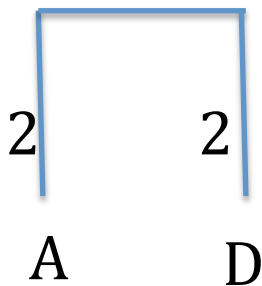
3.4 Add cluster  $C_k$  to the list of clusters, and remove clusters  $C_i$  and  $C_j$ .

Example: Consider the following distance matrix D:

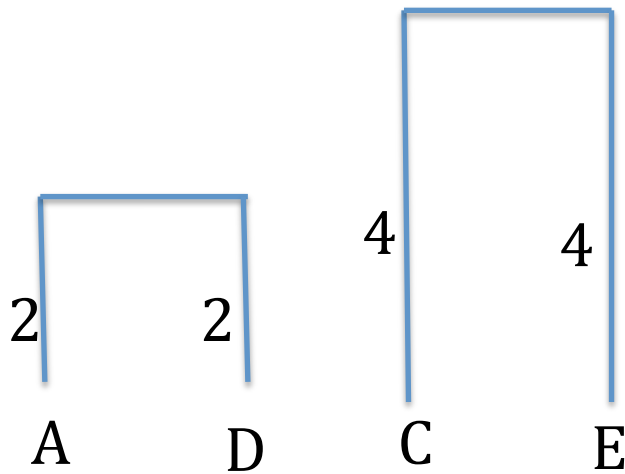
|   | A | B  | C  | D  | E  |
|---|---|----|----|----|----|
| A | - | 16 | 16 | 4  | 16 |
| B |   | -  | 10 | 16 | 10 |
| C |   |    | -  | 16 | 8  |
| D |   |    |    | -  | 16 |
| E |   |    |    |    | -  |

First set  $C_1=\{A\}$ ,  $C_2=\{B\}$ ,  $C_3=\{C\}$ ,  $C_4=\{D\}$ ,  $C_5=\{E\}$ .

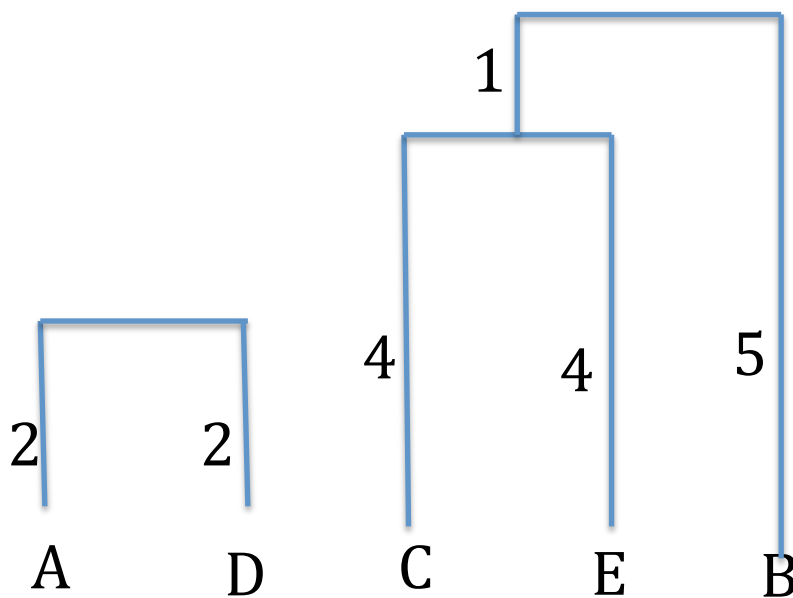
The pair with the smallest distance is  $(C_1, C_4)$ . Merge them to obtain  $C_6=\{A, D\}$  and create their parent node at distance  $d(C_1, C_4)/2 = 4/2 = 2$  from each, to obtain:



The next pair of clusters that is the closest is  $C_3$  and  $C_5$ , with  $d(C_3, C_5)=8$ . Merge them to obtain  $C_7=\{C, E\}$  and create their parent node at distance  $d(C_3, C_5)/2 = 8/2 = 4$  from each, to obtain:



The next pair of clusters that is the closest is  $C7=\{C,E\}$  and  $C2=\{B\}$ , with  $d(C7,C2)=(10+10)/2=10$ . Merge them to obtain  $C8=\{C,E,B\}$  and create their parent node at height  $d(C7,C2)/2 = 10/2 = 5$ , to obtain:



There are only two clusters  $C6$  and  $C8$ . Merge them and place their parent node at height  $d(C6,C8)/2 = ((16+16+16+16+16+16+16)/6) / 2 = 8$ , to obtain:

