

Distance-based methods for construction of phylogenetic trees

Lecture 6.4

by Marina Barsky

Defining distance

- How to measure distance between 2 DNA molecules so it reflects the time since they have evolved from a common ancestor?

The relative distance between genomes can be measured as the number of mutational events

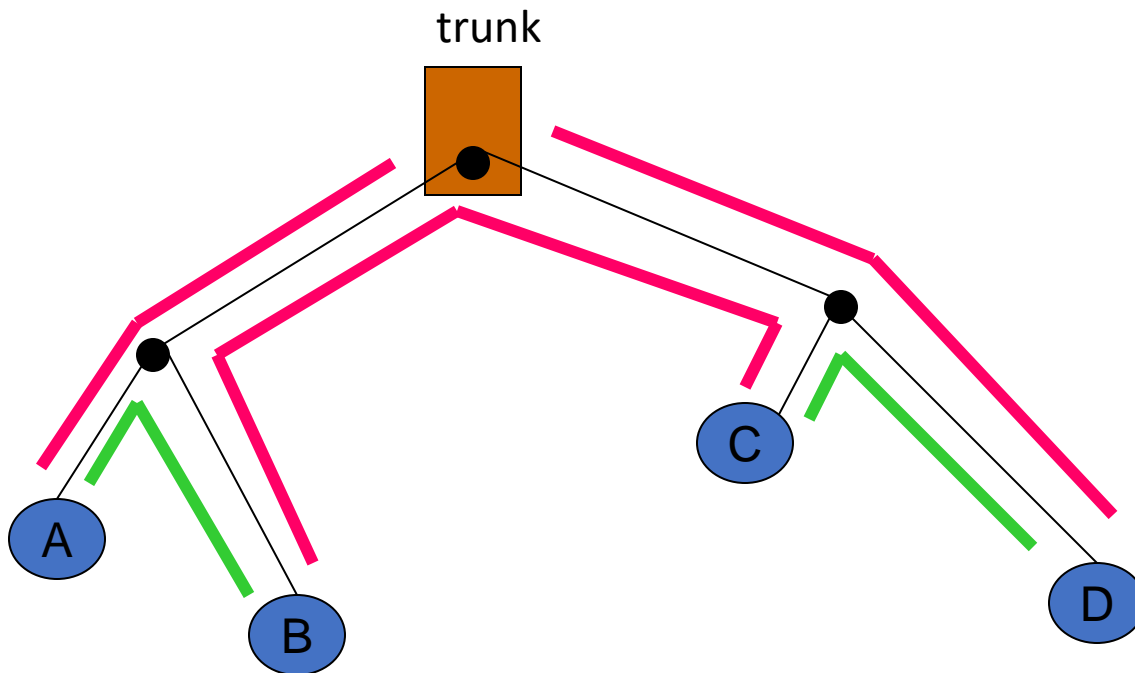
- Non computational: melting temperature of DNA hybrids
- Computational
 - Based on DNA or protein sequences
 - *Edit distance* – based on point mutations
 - Gene-sequence based
 - Alignment traces – align chromosomes from the different species, connect homologous genes by an edge. The number of crosses can be used as an evolutionary distance
- Better: combine different metrics

Building the tree

- Given a set of pairwise distances, find the best tree for a given data

Additive distances

- Distances that fit onto some tree are called *additive*. To determine if the distances are additive we use *the four point criterion*:

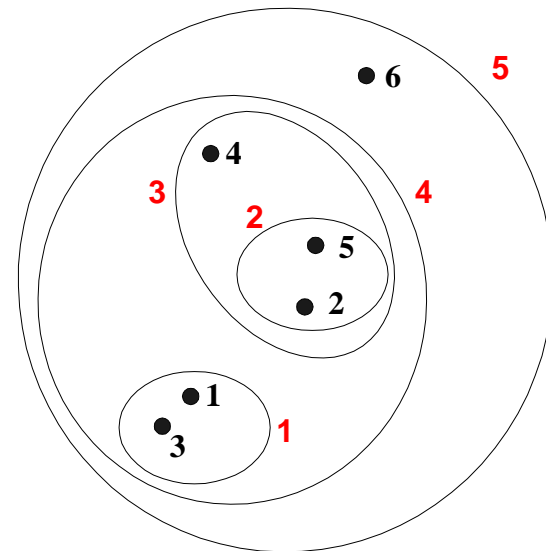
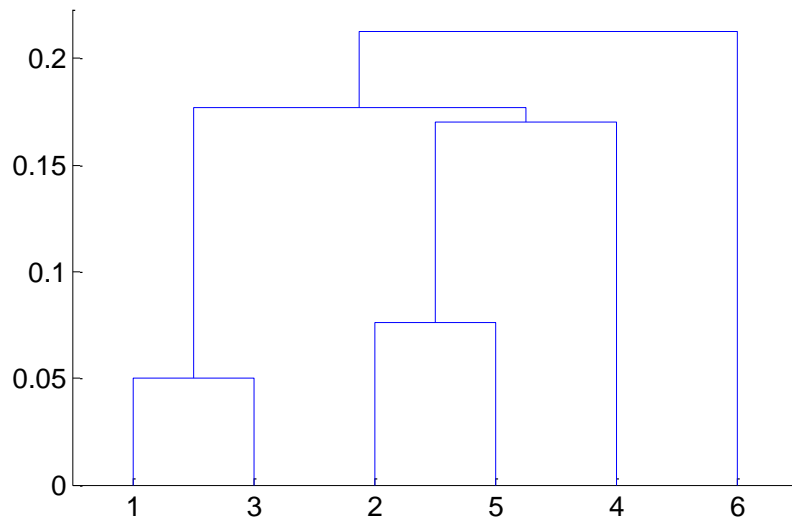


- The sums of pairwise distances that traverse the trunk are equal
- The sum of distances that traverse the trunk is \geq the sum of remaining distances

$$d_{AD} + d_{BC} = d_{BD} + d_{AC} \geq d_{AB} + d_{CD}$$

Hierarchical Clustering

- Produces a set of nested clusters organized as a hierarchical tree
- Can be visualized as a dendrogram
 - A tree-like diagram that records the sequences of merges or splits



Hierarchical Clustering

- Start with the points as individual clusters.
- At each step, merge the closest pair of clusters until only one cluster left.

Algorithm

Let each data point be a cluster

Compute the distance matrix

Repeat

Merge the two closest clusters

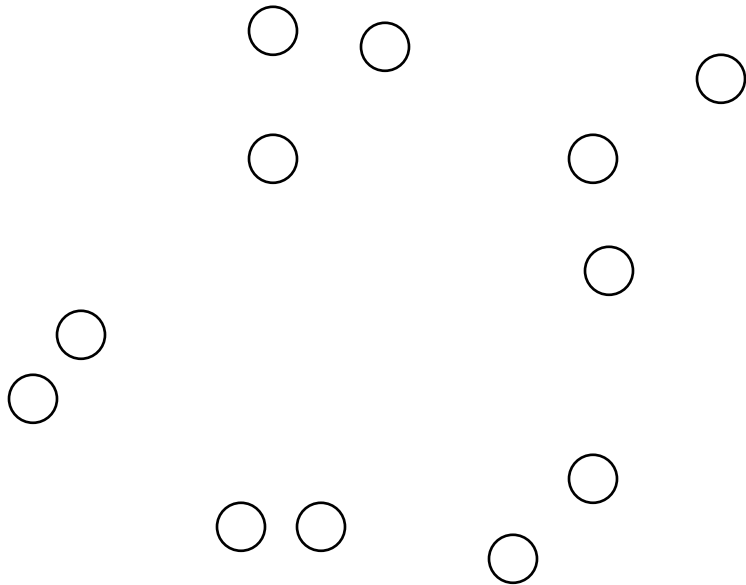
Update the distance matrix

Until only a single cluster remains

- Key operation is *the computation of the distance between two clusters.*

Starting Situation

- Start with clusters of individual points and a proximity matrix



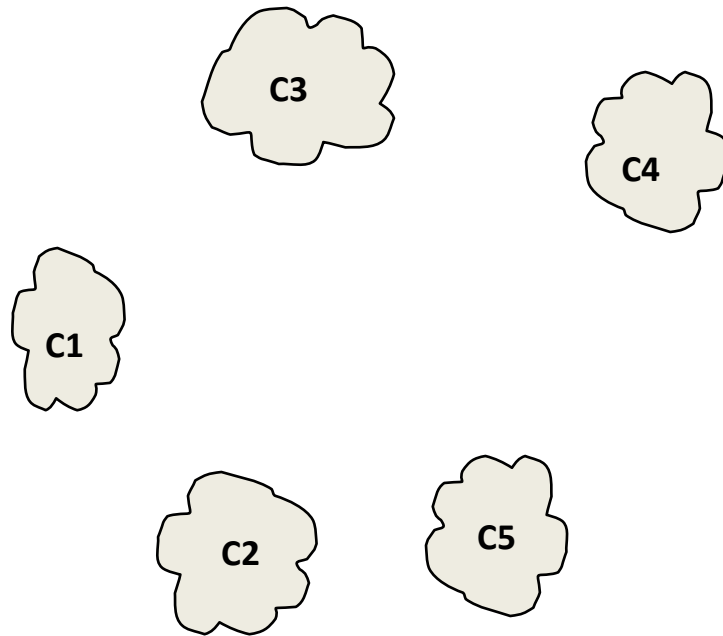
	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						
.						

Distance Matrix



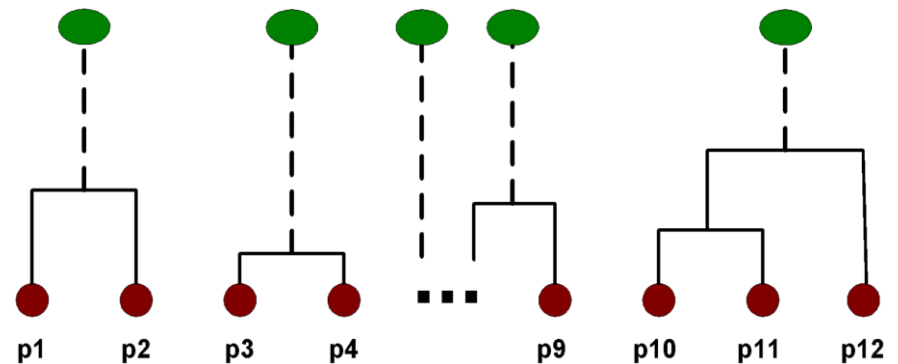
Intermediate Situation

- After some merging steps, we have some clusters



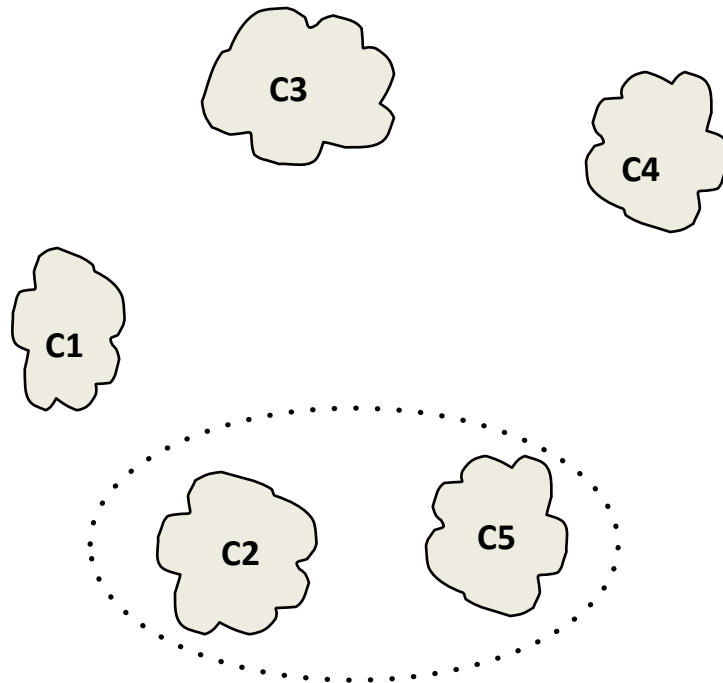
	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Distance Matrix



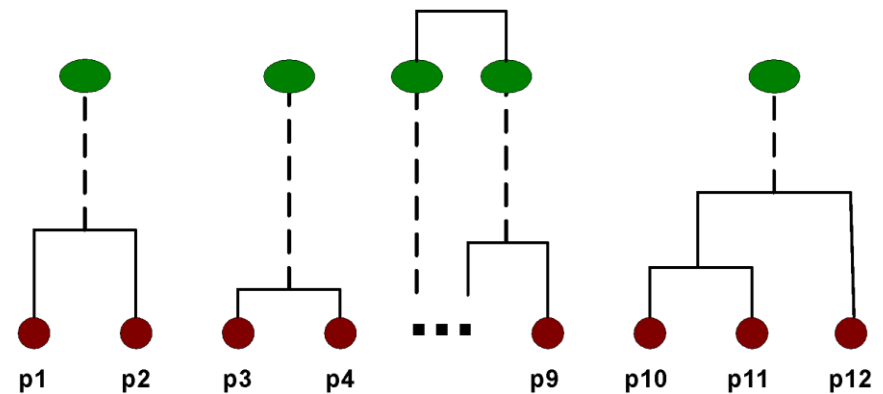
Intermediate Situation

- We want to merge the two closest clusters (C2 and C5) and update the distance matrix.



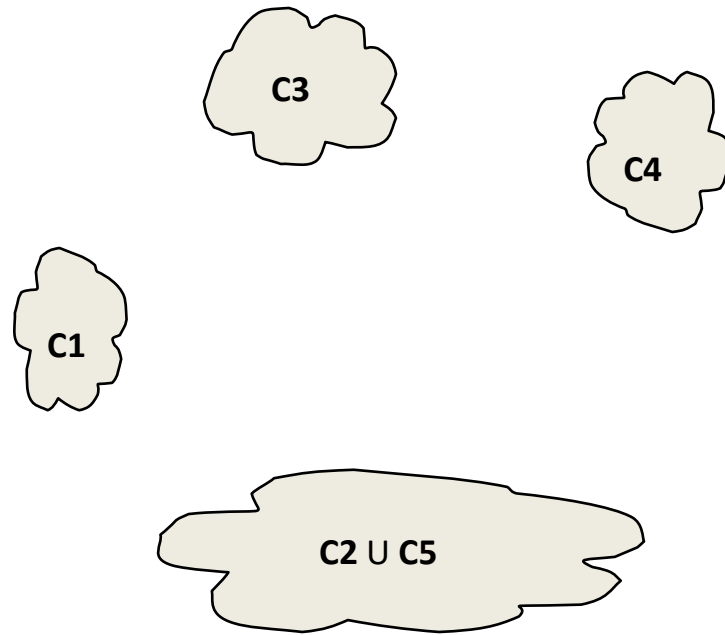
	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Distance Matrix



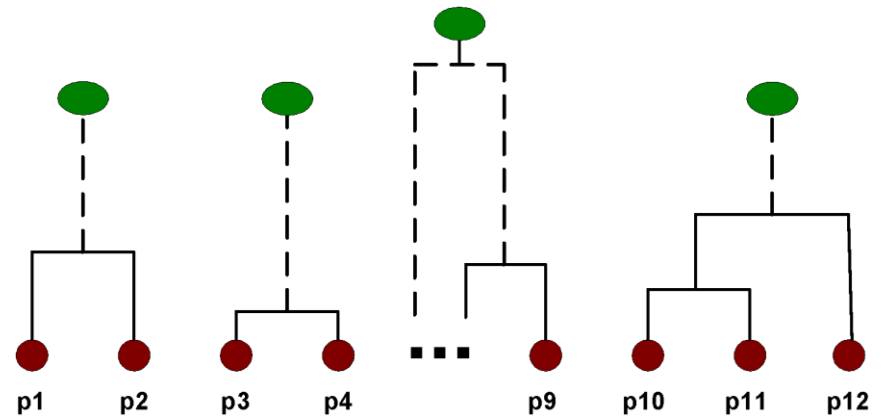
After Merging

- The question is “How do we update the distance matrix?”

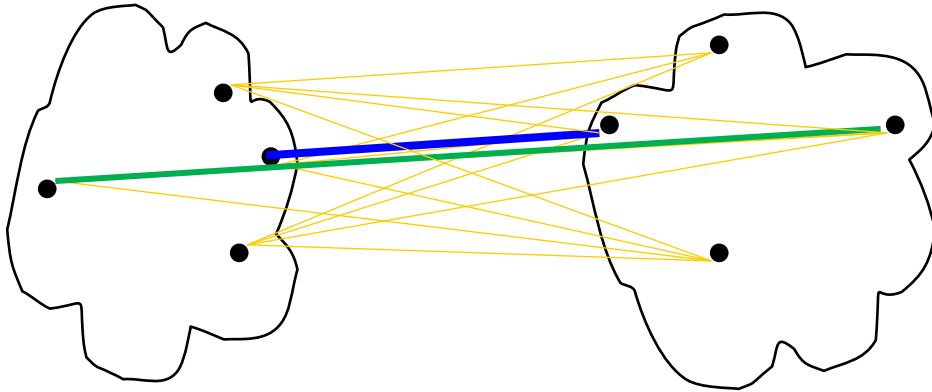


	C1	C2 U C5	C3	C4
C1		?		
C2 U C5	?		?	?
C3		?		
C4		?		

Distance Matrix



How to Define Inter-Cluster Distance



	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						

Distance Matrix

- **MIN**
- **MAX**
- **Group Average –UPGMA:**
Unweighted Pair-Group Method
using an arithmetic **Average**

Group Average

- Distance between two clusters is the average of all pairwise distances between points in the two clusters.

$$\text{distance}(\text{Cluster}_i, \text{Cluster}_j) = \frac{\sum_{\substack{p_i \in \text{Cluster}_i \\ p_j \in \text{Cluster}_j}} \text{distance}(p_i, p_j)}{|\text{Cluster}_i| * |\text{Cluster}_j|}$$

Example

	Human	Chimpanzee	Gorilla	Orangutan	Gibbon
Human					
Chimpanzee	1				
Gorilla	4	2			
Orangutan	8	7	5		
Gibbon	10	9	2	9	

The distances are real and determined based on the melting temperature of the DNA hybrids (the entire mitochondrial DNA)

Distance matrix

Distance matrix

	A	B	C	D	E
A					
B	1				
C	4	3			
D	8	7	2		
E	10	9	4	6	

Check: *Are the distances additive?*

Check each set of 4 points

ABCD

$$AC+BD=BC+AD=11 > AB+CD=3$$

ABCE

$$AC+BE=AE+BC=13 > AB+CE=5$$

ACDE

$$AD+CE=AE+CD=12 > AC+DE=10$$

BCDE

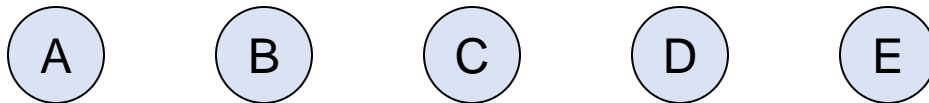
$$BD+CE=CD+BE=11 > BC+DE=9$$

	A	B	C	D	E
A					
B	1				
C	4	3			
D	8	7	2		
E	10	9	4	6	

**Yes, the distances are additive,
we can build the tree**

UPGMA – demo 1/4

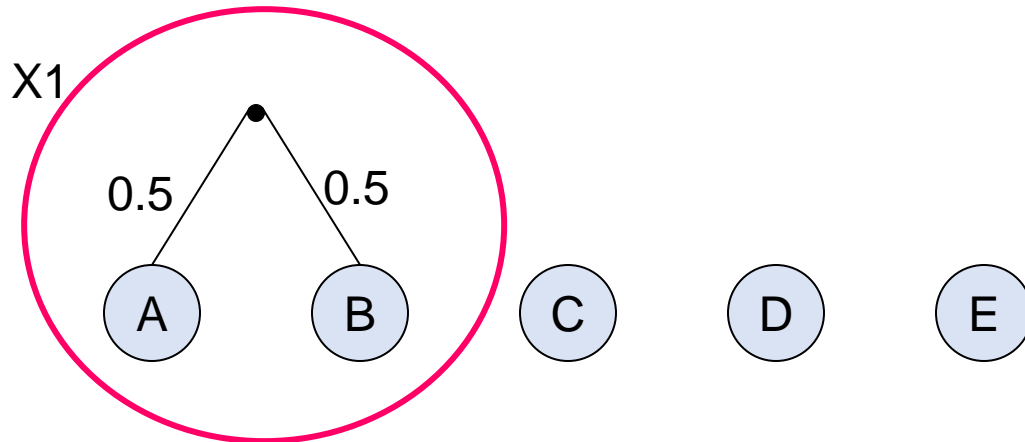
	A	B	C	D	E
A					
B	1				
C	4	3			
D	8	7	2		
E	10	9	4	6	



Basic clusters – distance 0 between the elements inside each cluster

UPGMA – demo 2/4

	A	B	C	D	E
A	0				
B	1	0			
C	4	3	0		
D	8	7	2	0	
E	10	9	4	6	0



Form cluster X1 with min distance between two points

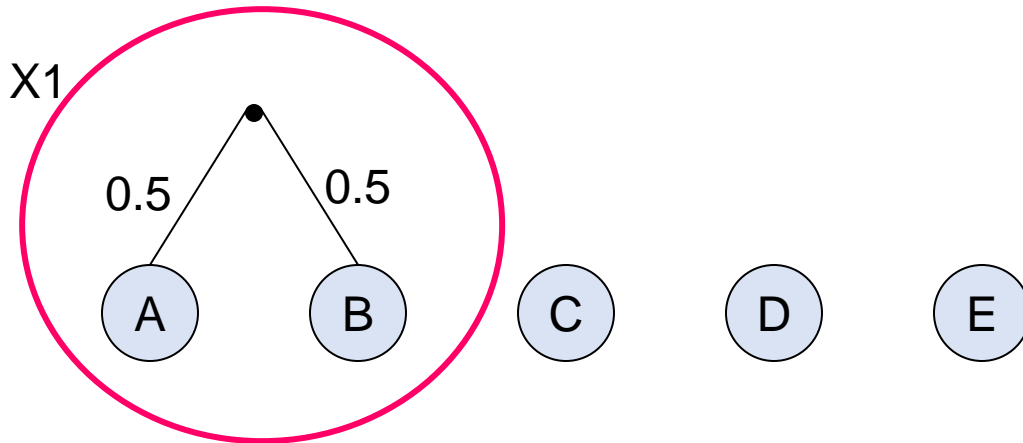
UPGMA – demo 2/4

	X1	C	D	E
X1	0			
C	3.5	0		
D	7.5	2	0	
E	9.5	4	6	0

Update distance matrix



	A	B	C	D	E
A	0				
B	1	0			
C	4	3	0		
D	8	7	2	0	
E	10	9	4	6	0



$$d_{X_1 C} = (d_{AC} + d_{BC}) / (2) * 1 = 3.5$$

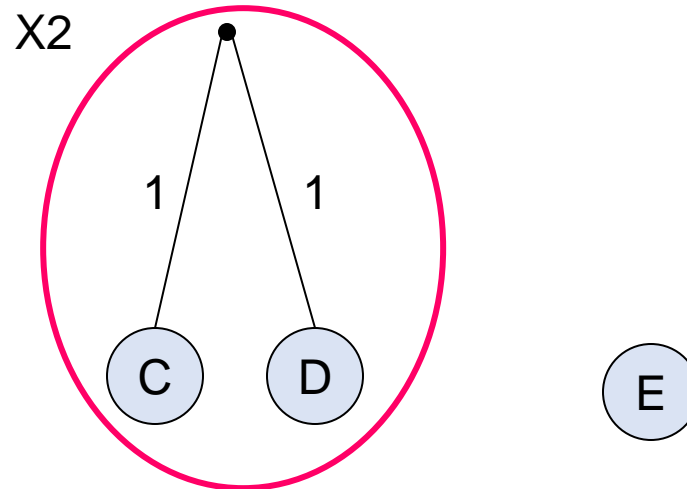
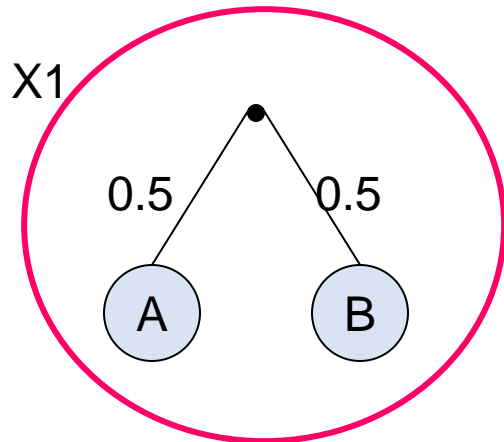
$$d_{X_1 D} = (d_{AD} + d_{BD}) / (2) * 1 = 7.5$$

$$d_{X_1 E} = (d_{AE} + d_{BE}) / (2) * 1 = 9.5$$

UPGMA – demo 3/4

	X1	C	D	E
X1	0			
C	3.5	0		
D	7.5	2	0	
E	9.5	4	6	0

	A	B	C	D	E
A	0				
B	1	0			
C	4	3	0		
D	8	7	2	0	
E	10	9	4	6	0



Form cluster X2 with min distance between two points

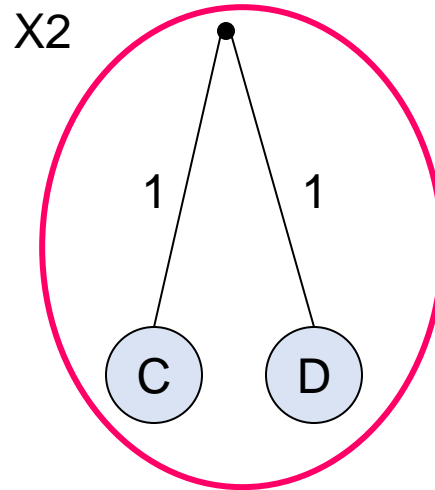
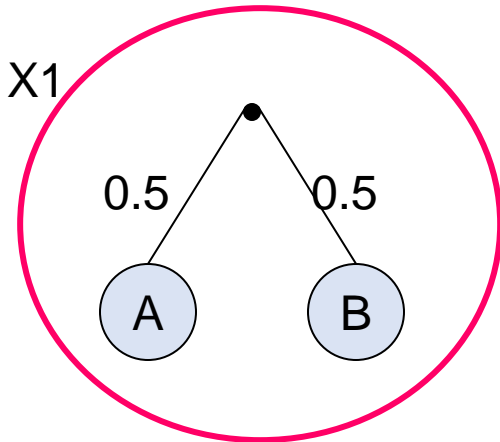
UPGMA – demo 3/4

	X1	X2	E
X1	0		
X2	5.5	0	
E	7.5	5	0

Update distance matrix

	X1	C	D	E
X1	0			
C	3.5	0		
D	7.5	2	0	
E	9.5	4	6	0

	A	B	C	D	E
A	0				
B	1	0			
C	4	3	0		
D	8	7	2	0	
E	10	9	4	6	0



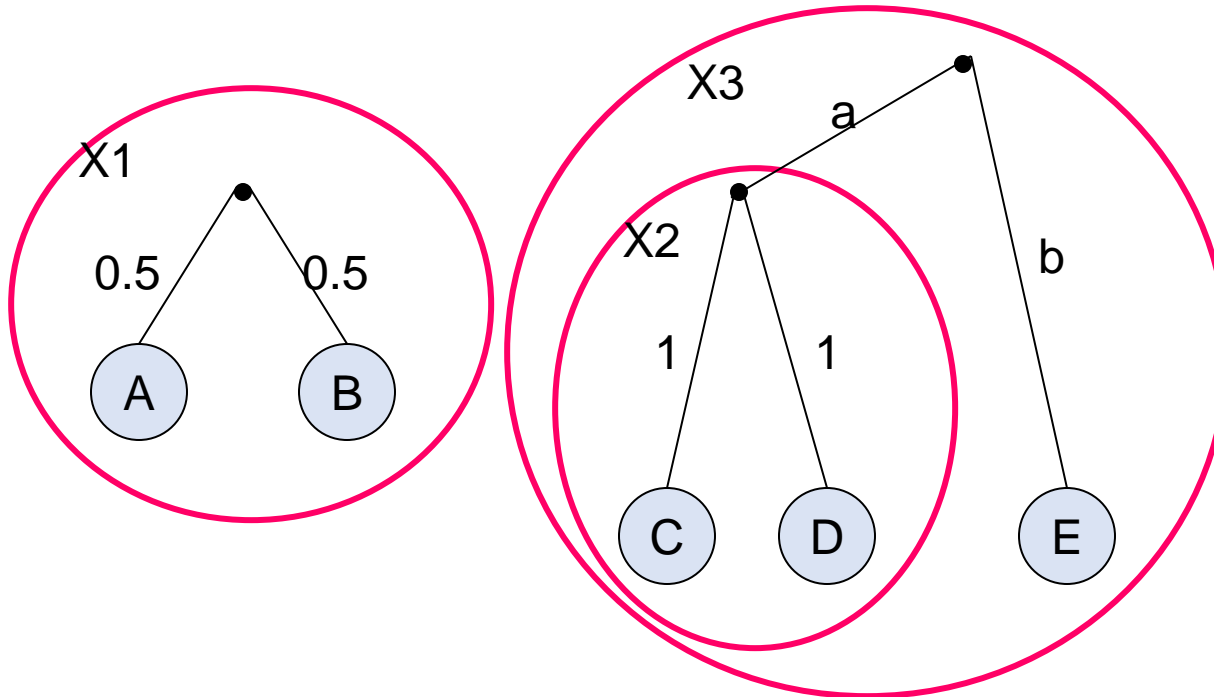
$$d_{X_2 E} = (d_{CE} + d_{DE}) / (2) * 1 = 5$$

$$d_{X_1 X_2} = (d_{AC} + d_{BC} + d_{AD} + d_{BD}) / (2 * 2) = 5.5$$

UPGMA – demo 4/4

	X1	X2	E
X1	0		
X2	5.5	0	
E	7.5	5	0

	A	B	C	D	E
A	0				
B	1	0			
C	4	3	0		
D	8	7	2	0	
E	10	9	4	6	0



$$a = 1/2 * 5 - 1 = 1.5$$

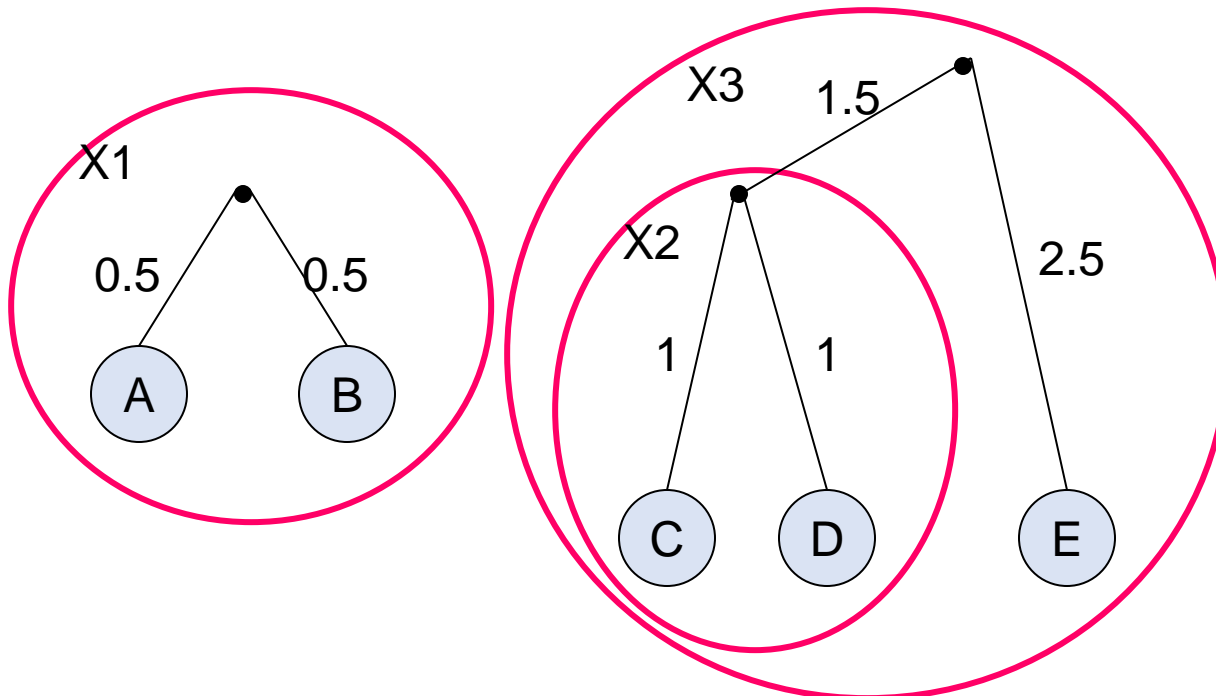
$$b = 1/2 * 5 = 2.5$$

Create cluster X3 with min distance between 2 clusters

UPGMA – demo 4/4

	X1	X2	E
X1	0		
X2	5.5	0	
E	7.5	5	0

	A	B	C	D	E
A	0				
B	1	0			
C	4	3	0		
D	8	7	2	0	
E	10	9	4	6	0



Distribute the distance between 2 edges to have half of this distance from the root to leaves in both branches

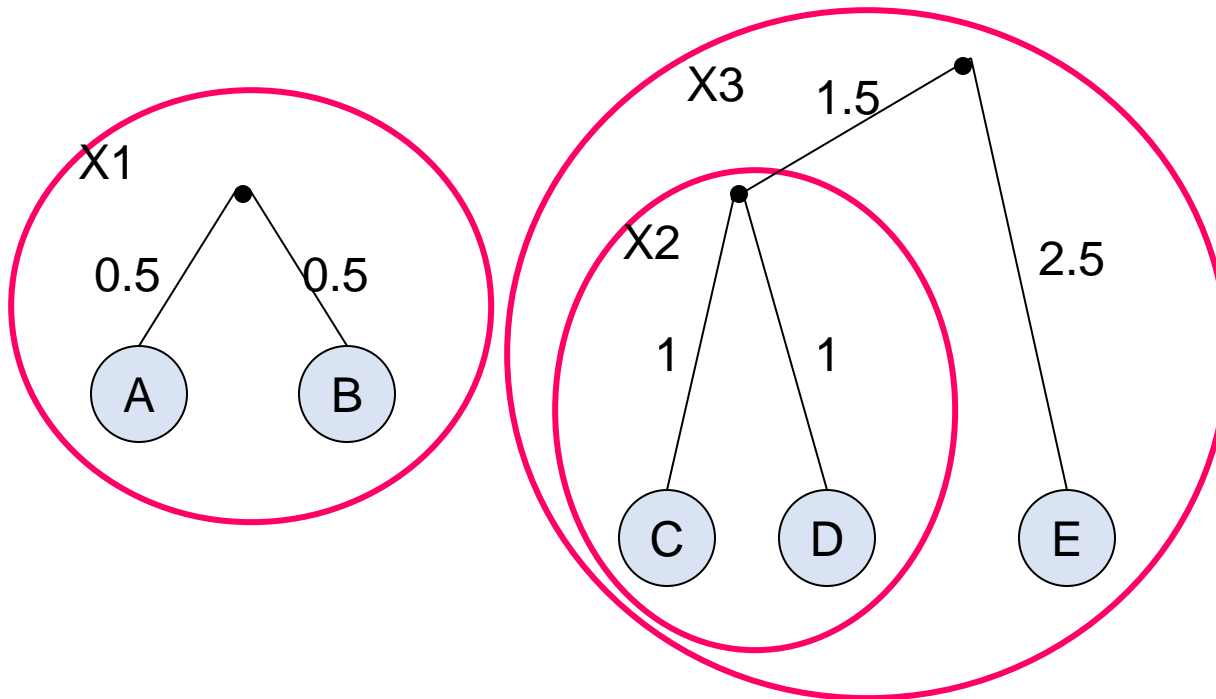
UPGMA – demo 4/4

	X1	X3
X1	0	
X3	6.8	0



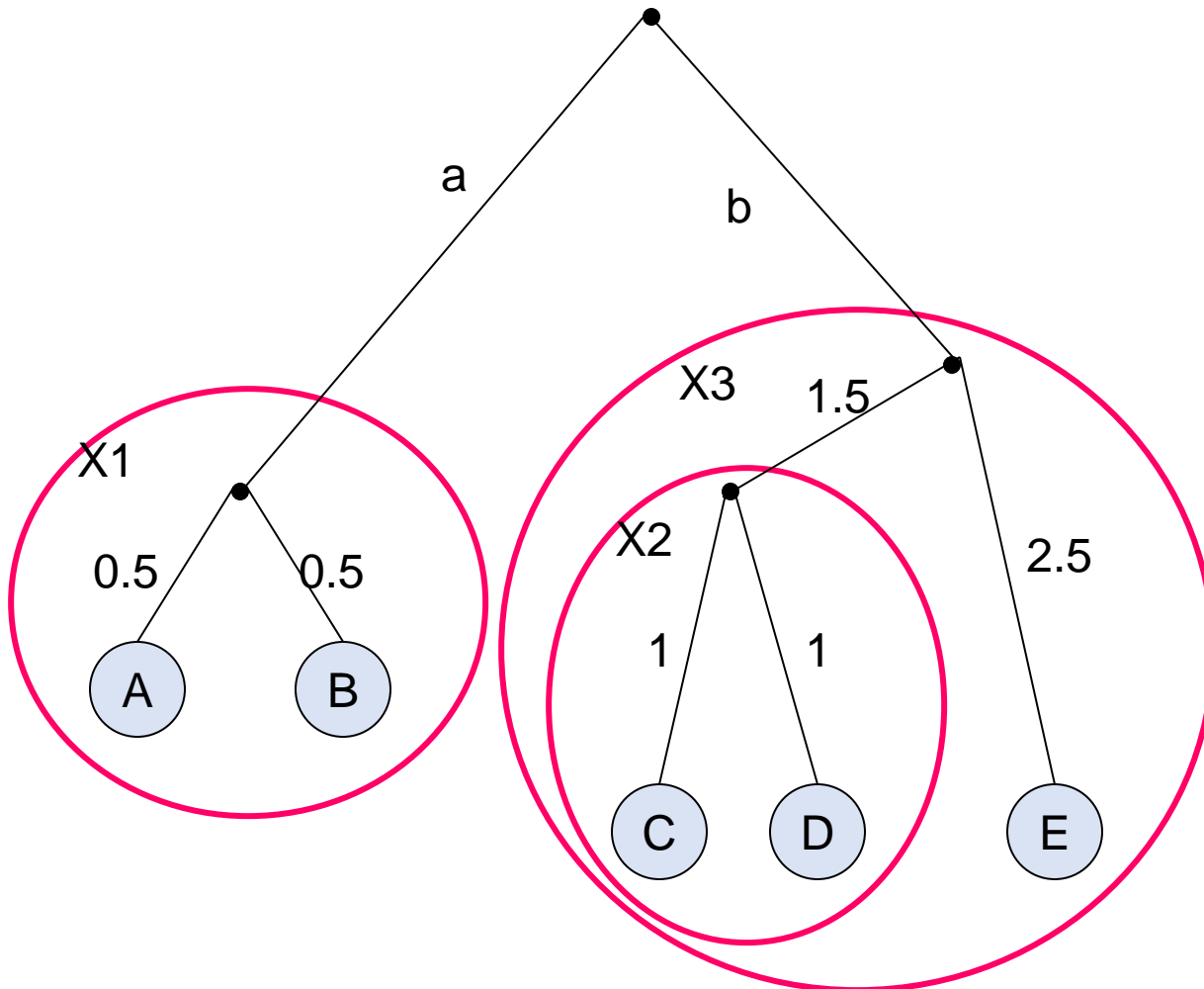
	A	B	C	D	E
A	0				
B	1	0			
C	4	3	0		
D	8	7	2	0	
E	10	9	4	6	0

Update distance matrix



$$d_{X1X3} = (d_{AC} + d_{AD} + d_{AE} + d_{BC} + d_{BD} + d_{BE}) / 2 * 3 = (4 + 8 + 10 + 3 + 7 + 9) / 6 = 6.8$$

UPGMA – demo 4/4



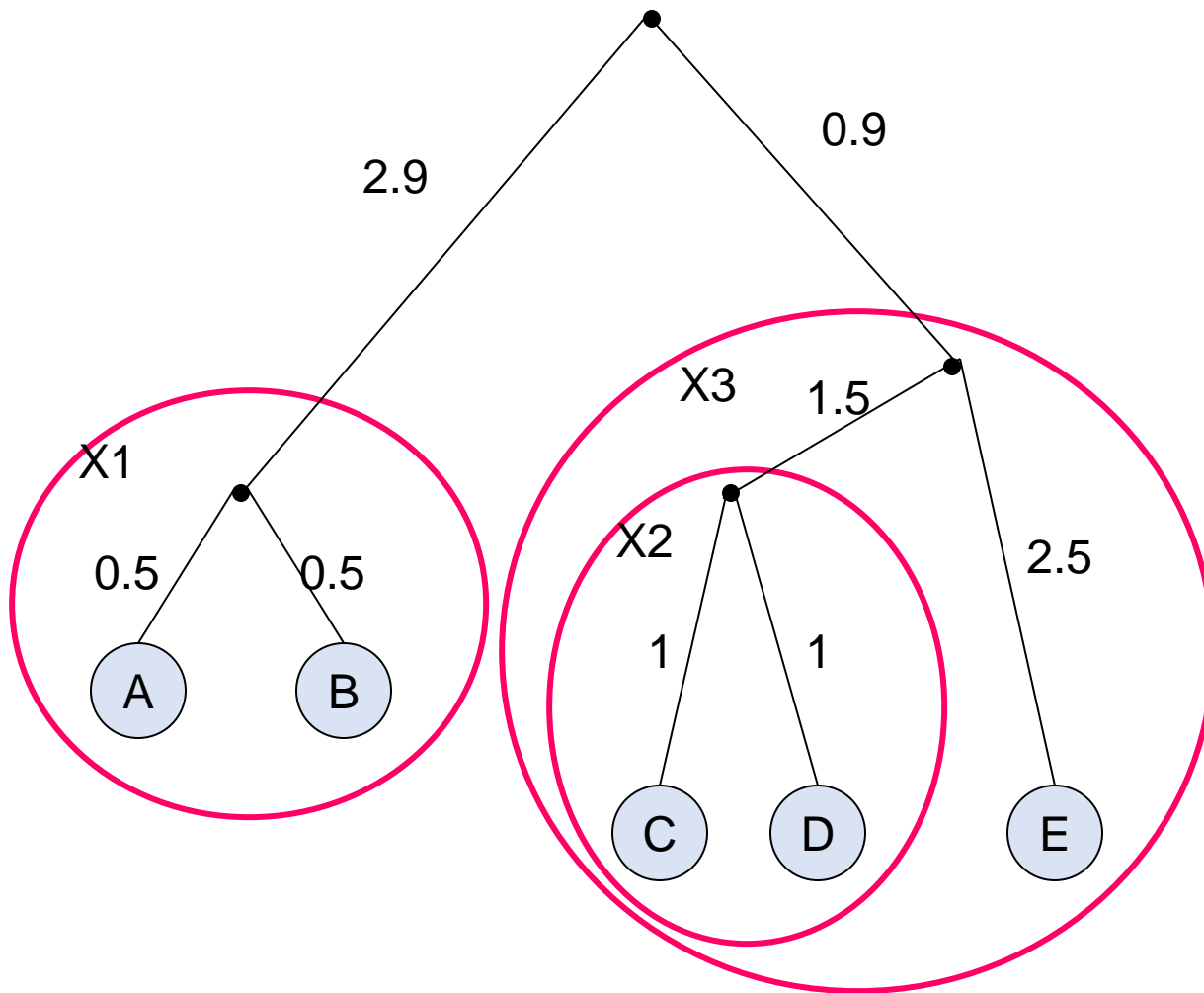
	X1	X3
X1	0	
X3	6.8	0

$$a = 6.8/2 - 0.5 = 2.9$$

$$b = 6.8/2 - 2.5 = 0.9$$

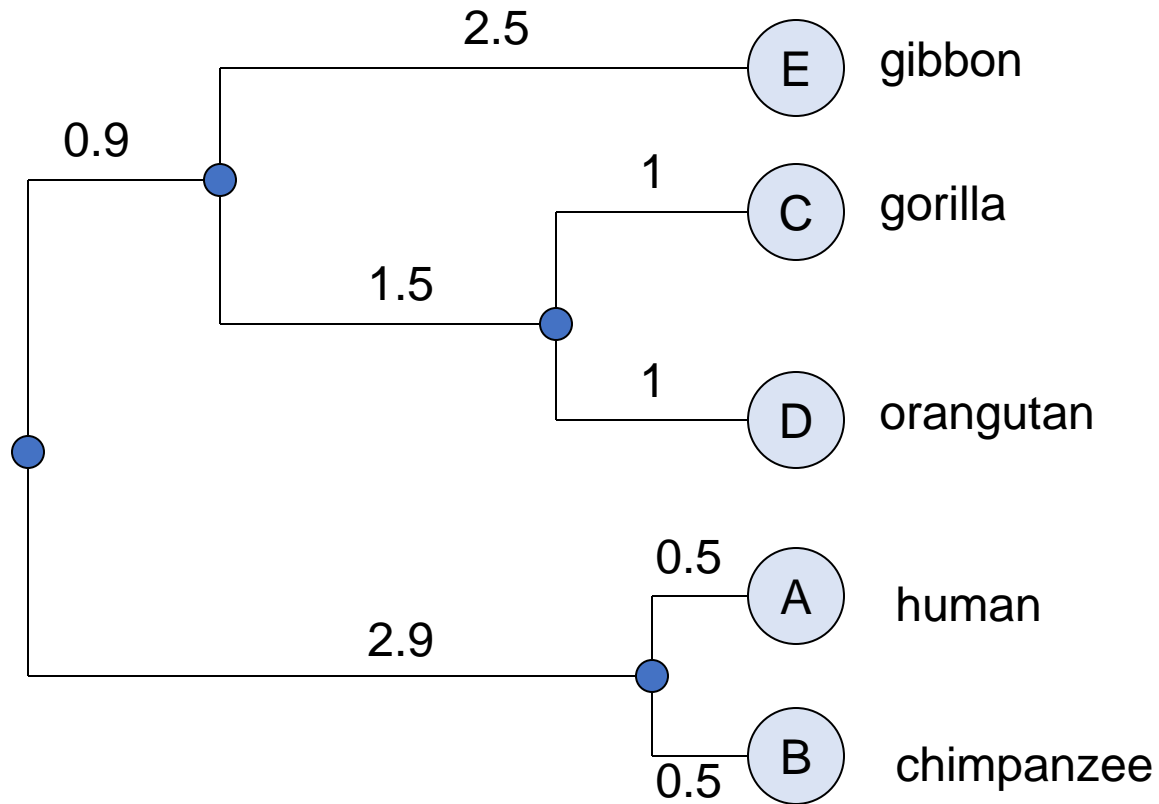
Distribute the distance between 2 edges to have half of this distance from the root to leaves in both branches

UPGMA – the resulting tree

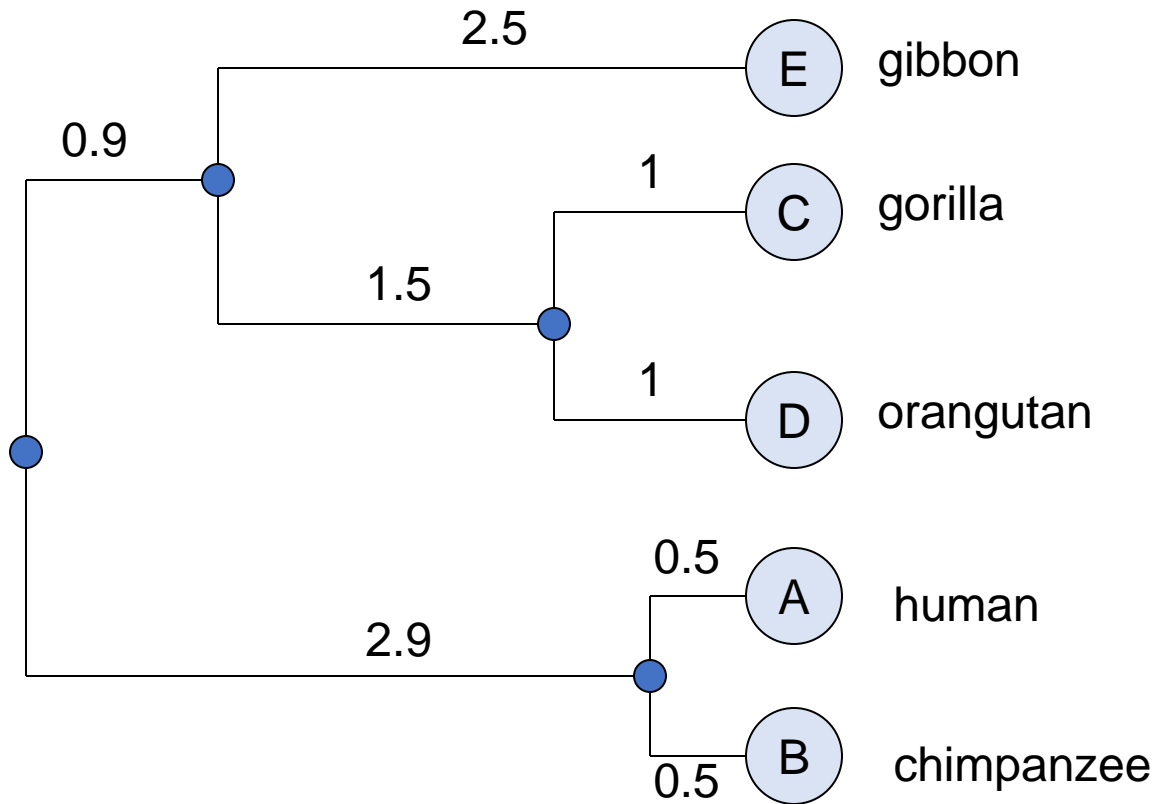


	A	B	C	D	E
A	0				
B	1	0			
C	4	3	0		
D	8	7	2	0	
E	10	9	4	6	0

The resulting tree with distances



Molecular clock



The edge lengths can be viewed as times measured by *molecular clock* with a constant rate of mutational events.

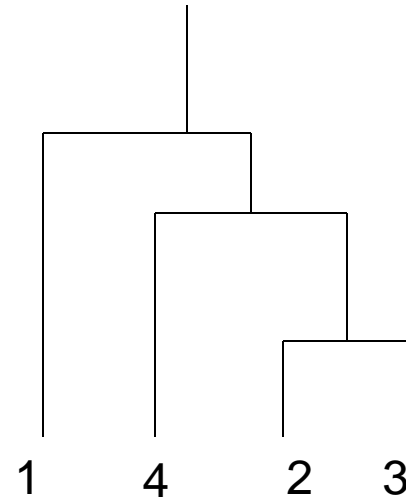
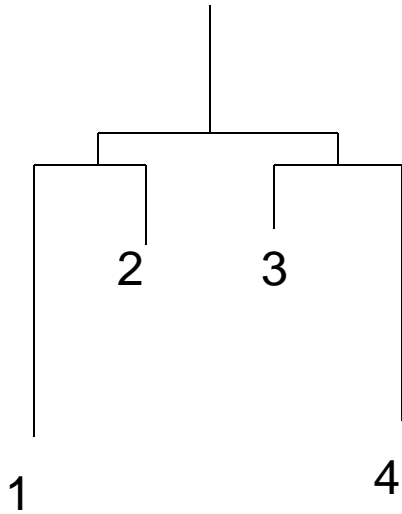
We assume that divergence occurred at the same time at all branching points, the sum of edge lengths from any node to all leaves is the same for any possible path

Hierarchical Clustering: Time and Space

- $O(N^2)$ space since it uses the proximity matrix.
 - N is the number of data points.
- $O(N^3)$ time in many cases
 - There are N steps and at each step the size, N^2 , proximity matrix must be updated and searched
 - Complexity can be reduced to $O(N^2 \log(N))$ time using more advanced data structures

Hierarchical clustering algorithm is expensive !

When the tree does not reflect reality



- If the original tree, which we try to reconstruct, had different path lengths to its leaves, it will be reconstructed incorrectly by UPGMA.
- In this case, the closest leaves (2,3) are not siblings and they do not have a common parent, which will be assigned to them by UPGMA

Predict correctness:

Test for ultrametric condition

- We can predict whether the reconstruction of the real tree is likely to be correct by **testing** our distances **for ultrametric condition**:

The distance matrix is *ultrametric* if for **any triplet** of sequences, X_i, X_j, X_k , the distances d_{ij}, d_{ik}, d_{jk} are either all equal or two are equal and the remaining one is smaller

Thus, if distances were derived from a real tree with a molecular clock, the distance matrix will be ultrametric

Ultrametric and non-ultrametric distance matrices

	A	B	C	D
A	0			
B	1	0		
C	4	2	0	
D	8	7	5	0

$d_{AB}=1$, $d_{AC}=4$, $d_{BC}=2$

Non-ultrametric matrix ❌

	A	B	C	D
A	0			
B	4	0		
C	2	4	0	
D	8	8	8	0

$d_{AB}=4$, $d_{AC}=2$, $d_{BC}=4$ ✔️

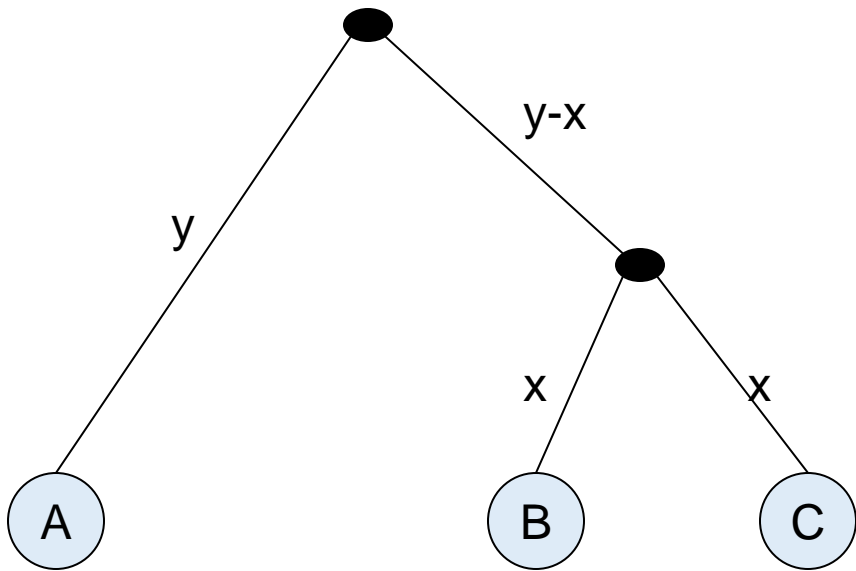
$d_{AC}=2$, $d_{AD}=8$, $d_{CD}=8$ ✔️

$d_{BC}=4$, $d_{BD}=8$, $d_{CD}=8$ ✔️

Ultrametric matrix

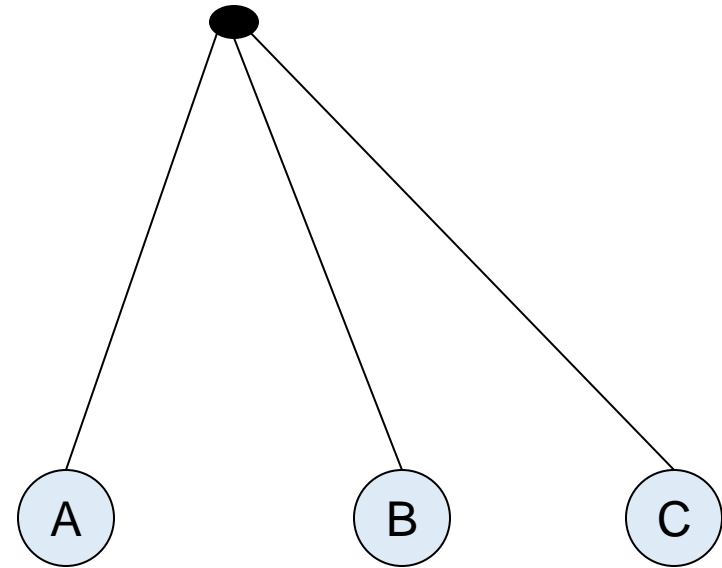
Ultrametric trees: why

$$AB=AC>BC$$



$$AB=AC=BC$$

or



$$AB=y+y-x+x=2y$$

$$AC=2y$$

$$BC=2x, x \leq y, \text{ since } y-x \geq 0 \text{ (no negative edge lengths)}$$

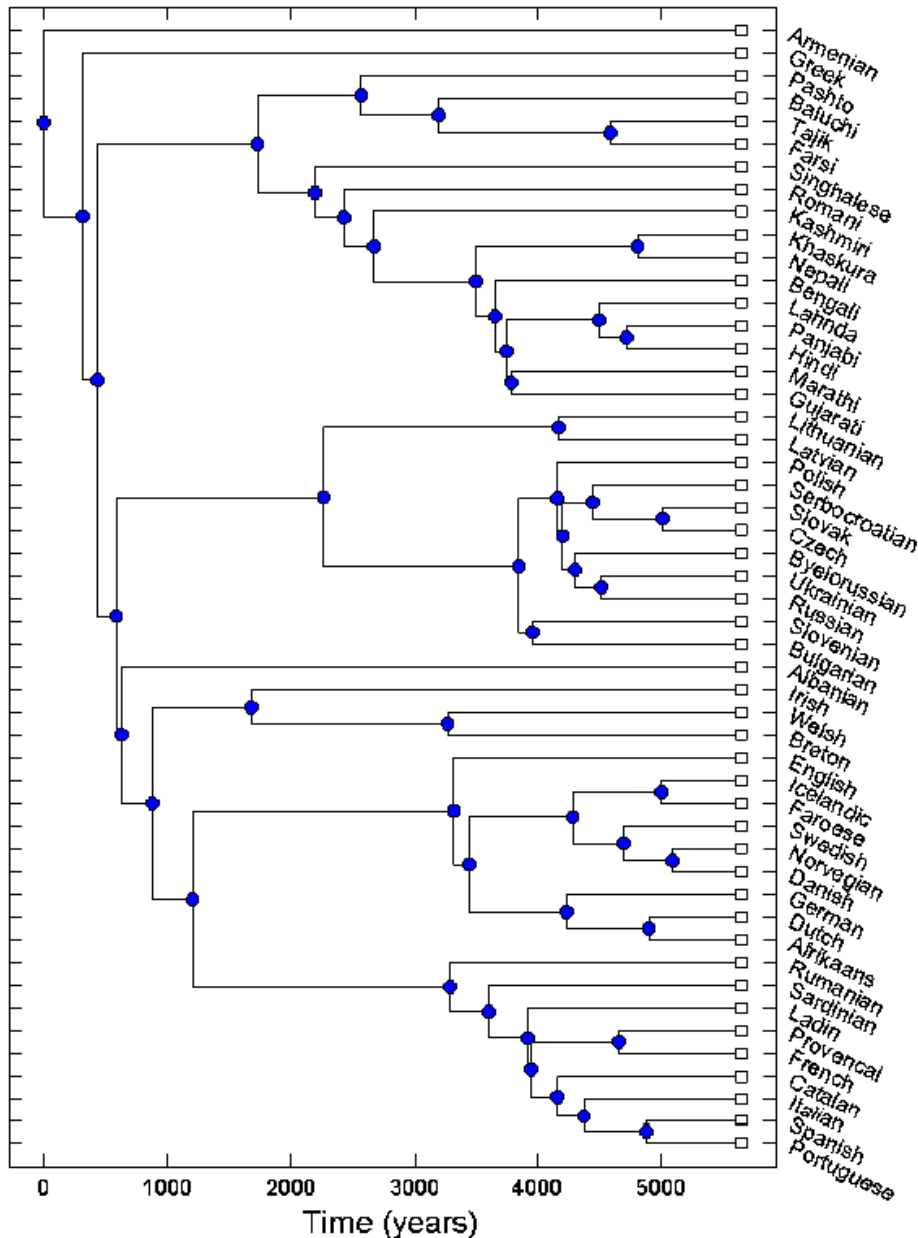
The rule for ultrametric trees:

2 out of 3 distances have a tie, and are \geq than the third distance

What to do if the distances are not ultrametric?

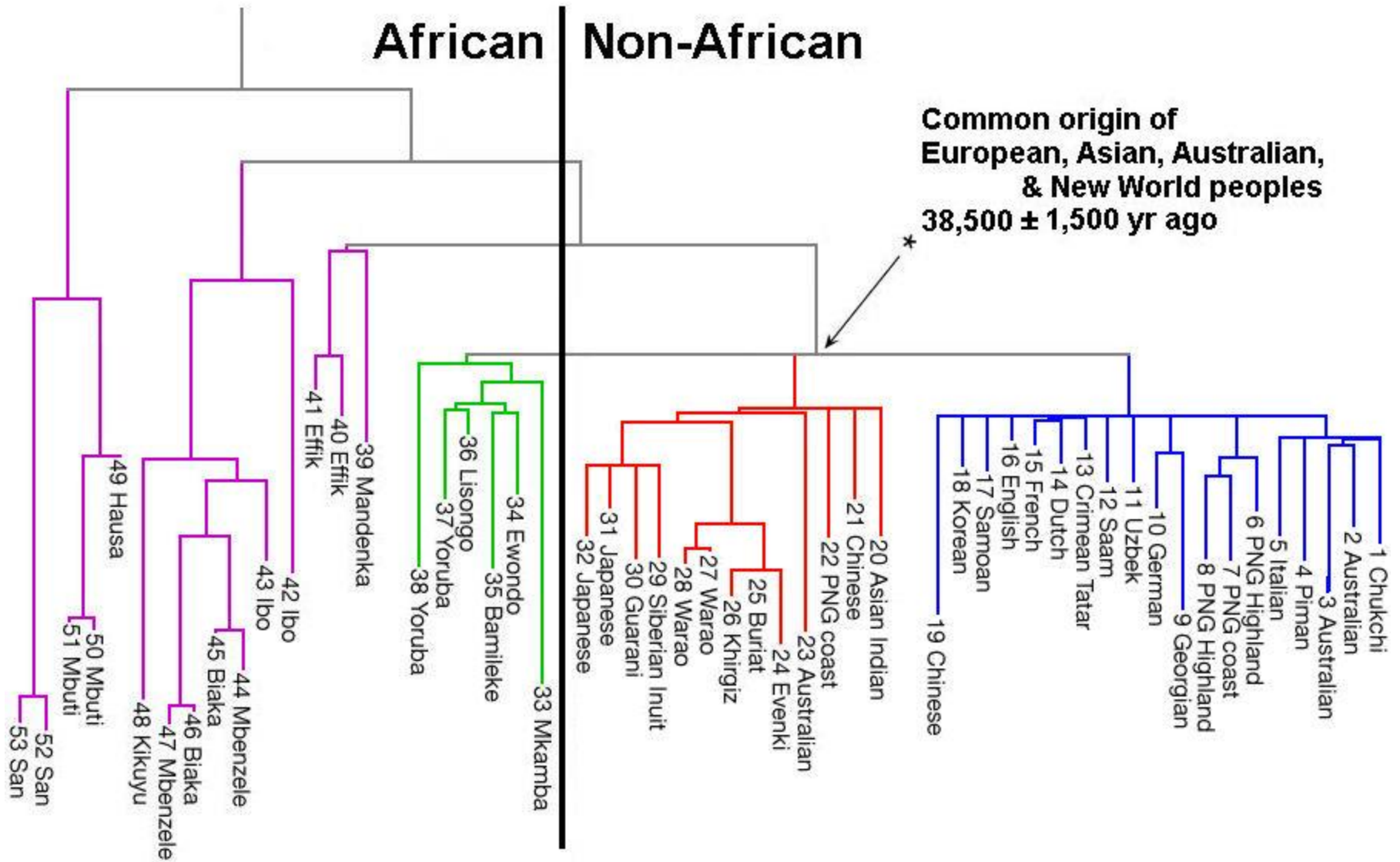
- We can apply UPGMA, but it may produce an incorrect tree.
- **The neighbor-joining algorithm** produces better results for non-ultrametric distances (see textbook)
- The main feature of the neighbor-joining algorithm is that we take into account not only how close are two clusters, but also how far away are they from other clusters
- This algorithm produces unrooted trees

Evolution of languages

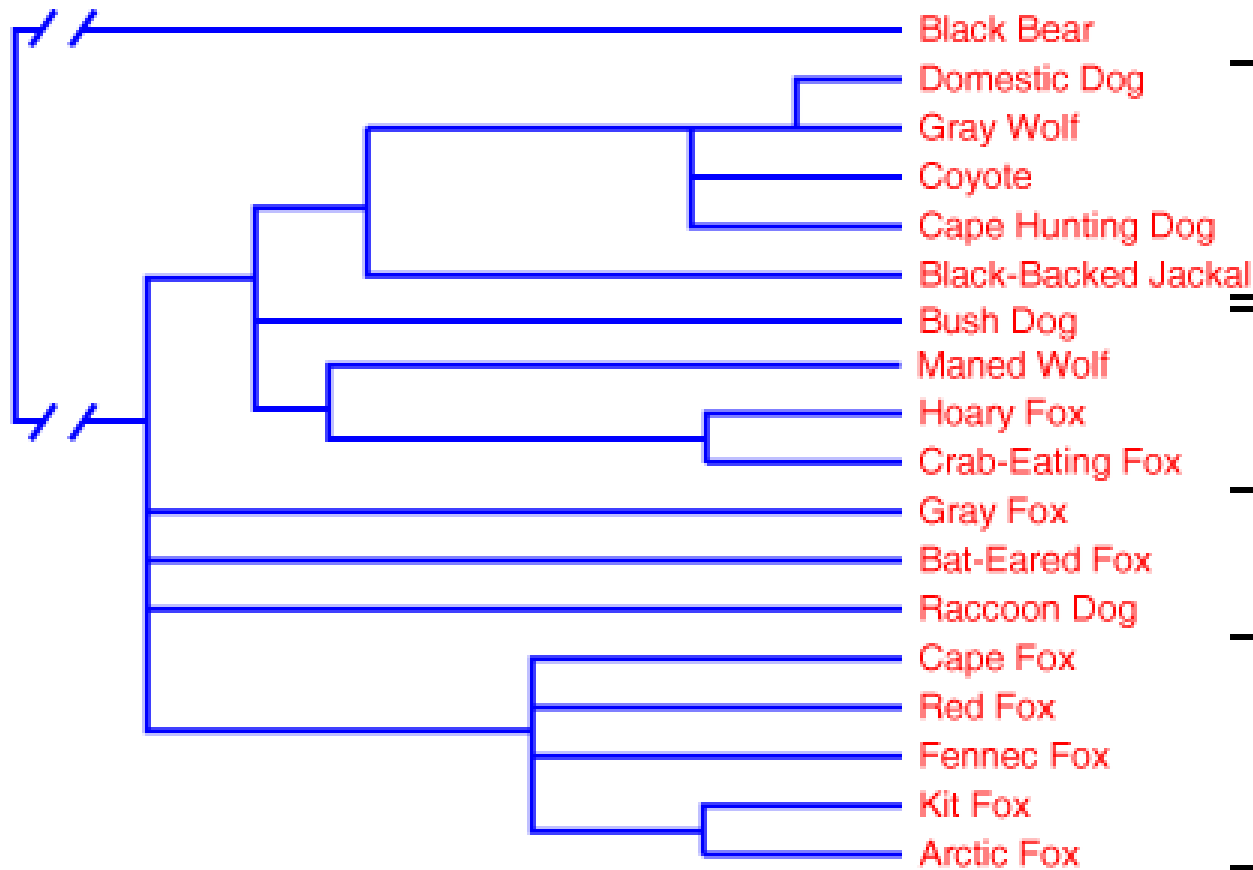


From
“Indo-European
languages tree by
Levenshtein distance”
by M. Serva1 and F.
Petroni

Evolution of human populations



Evolution of Canidae



Wolf-like canids



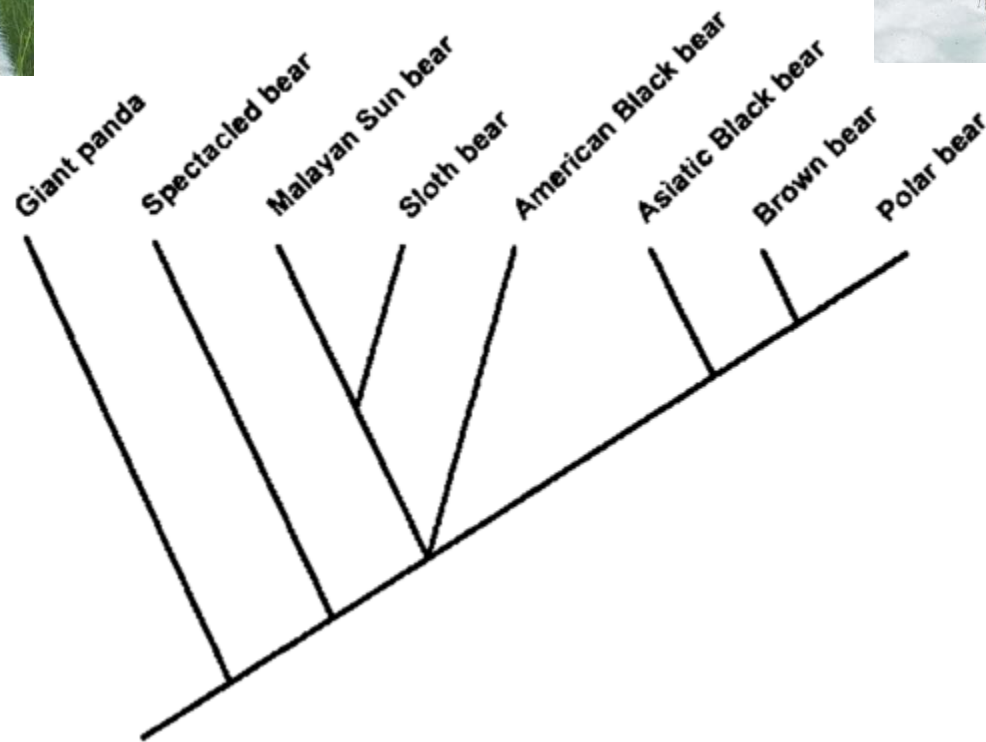
South American canids



Fox-like canids



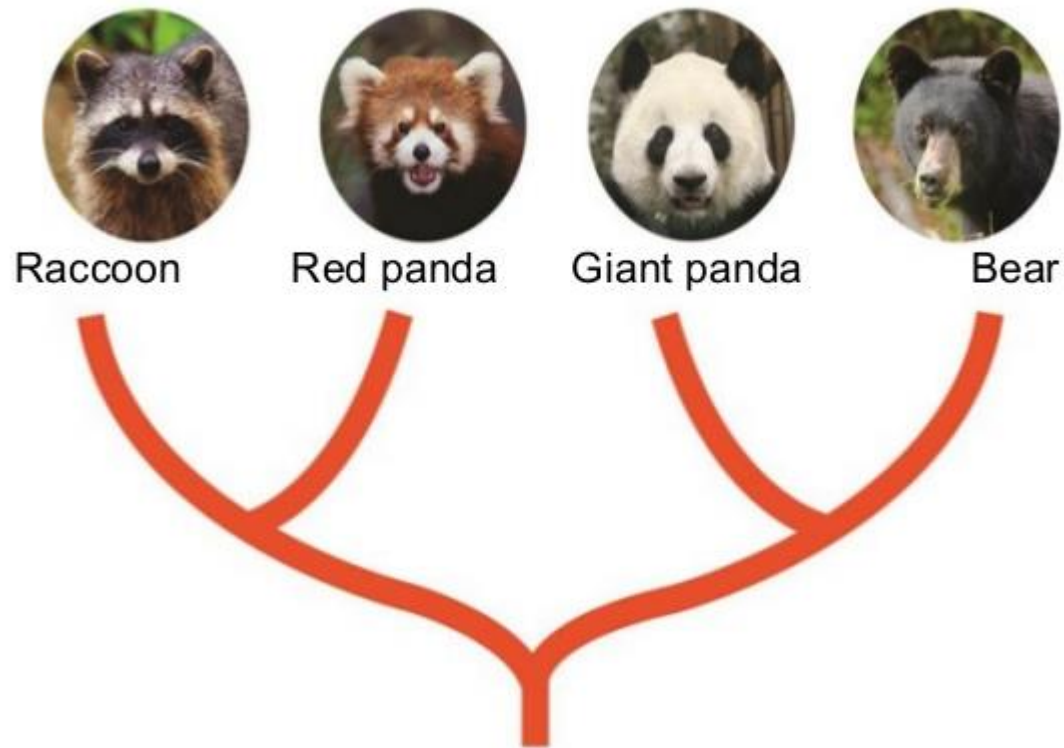
Giant Panda is a bear



What about Red Panda: a Cat or a Bear?



Red Panda: a Bear or a Cat?



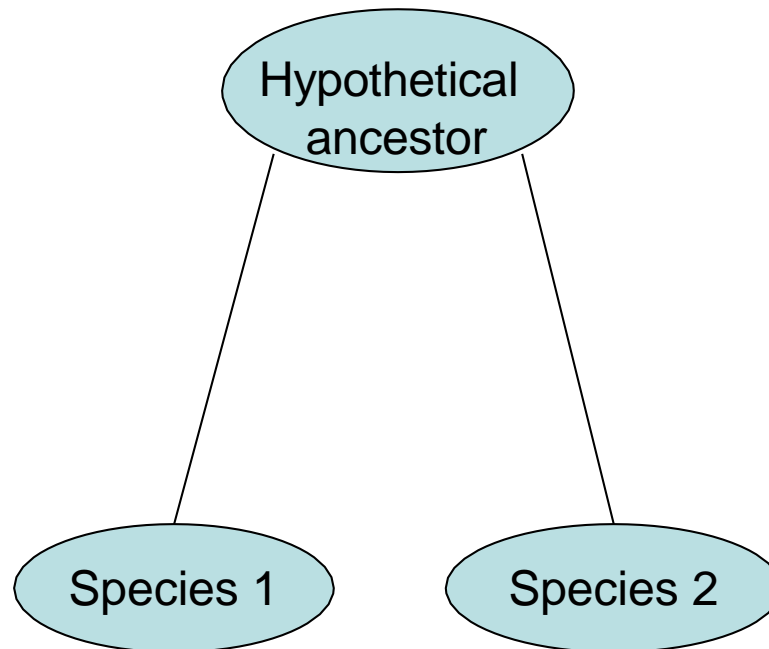
Flynn, J. J.; Nedbal, M. A.; Dragoo, J. W.; Honeycutt, R. L. (2000). "Whence the Red Panda?" Molecular Phylogenetics and Evolution.

Phylogenetic trees

Summary

We are solving the following problem:

- Inferring phylogenies from a given data set



Two main approaches to inferring phylogenies

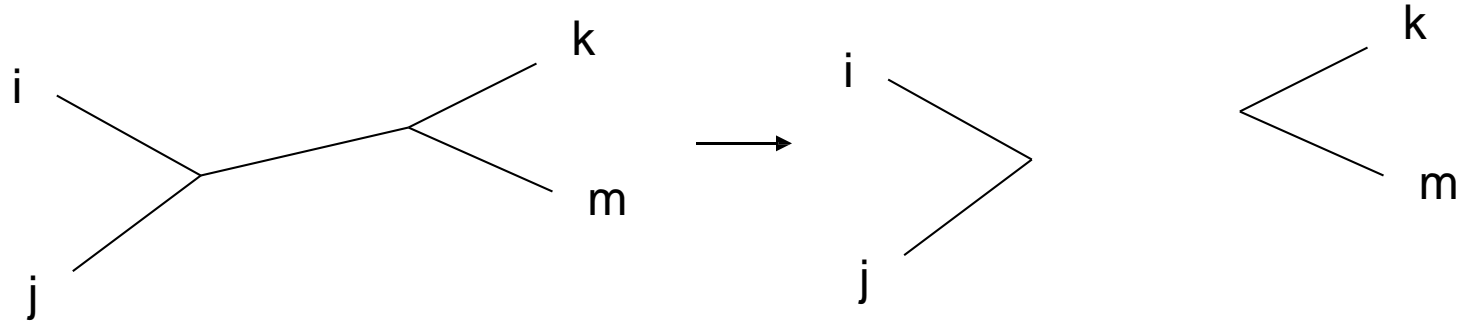
- **Distance-based**: for example pairwise edit distance, score of pairwise alignment etc.
- **Character-based**: examine each character separately in a given site of a biological sequence

Distance-based - recap

- Input: distance matrix of pairwise distances for N species
- Goal: find a tree consistent with the distance matrix. This means that the sum of edge lengths connecting each pair of leaves ij corresponds to a distance M_{ij}

Additivity of distances

- For any 4 objects, i, j, k, m , there are 3 different pairwise sums:
 $D_{ij} + D_{km}$, $D_{ik} + D_{jm}$, $D_{im} + D_{jk}$
- From these 3 sums, 2 should be equal and greater than the third – this will allow to group the smallest in between the 3 into a separate subtree

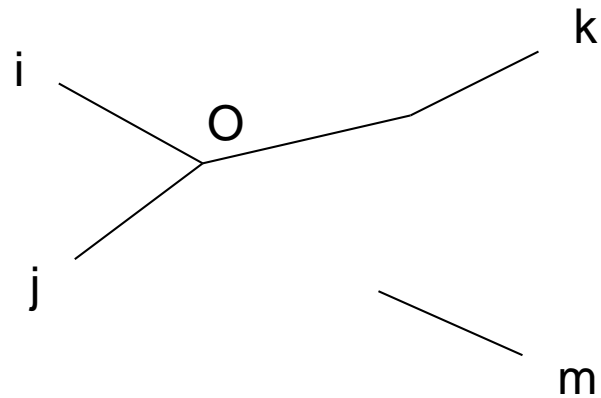


The distances are additive iff for any 4 objects there exists the following combination: $D_{ij} + D_{km} < (D_{im} + D_{jk} = D_{ik} + D_{jm})$

Why the distances have to be additive

For 3 objects, any set of distances is OK:

Let $iO=a$, $jO=b$, $kO=c$, then $D_{ij}=a+b$, $D_{ik}=a+c$, $D_{jk}=b+c$

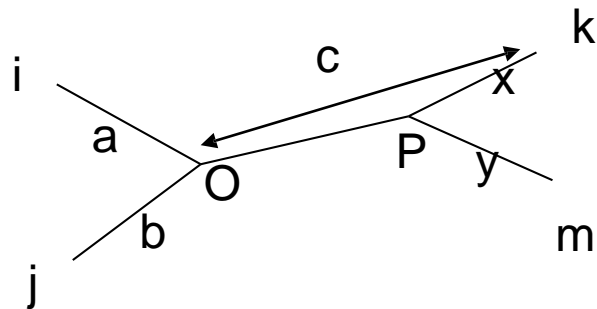


Why the distances have to be additive

Let $iO=a$, $jO=b$, $kO=c$

$D_{ij}=a+b$, $D_{ik}=a+c$, $D_{jk}=b+c$

Let's add the fourth object, m , to an arbitrary point P in the tree, and let $mP=y$, and $kP=x$, then $OP=c-x$



Then:

$$D_{im}+D_{jk}=a+c-x+y+b+c=a+b+y+(2c-x)$$

$$D_{ik}+D_{jm}=a+c+b+c-x+y=a+b+y+(2c-x)$$

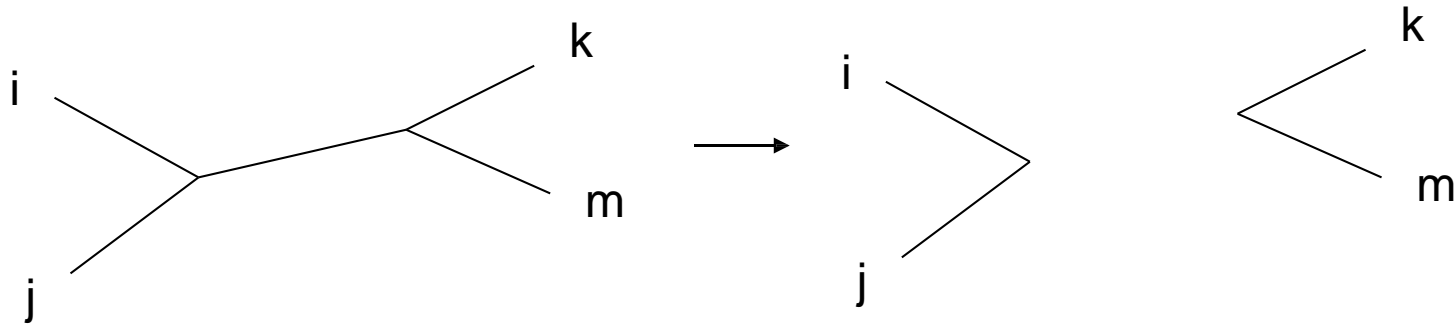
and

$$D_{ij}+D_{km}=a+b+x+y=a+b+y+x$$

$$2c-x \geq x, \text{ since } 2c \geq 2x$$

The rule of additivity

- The distances are additive if, for any 4 objects, i, j, k, m ,



$$D_{ij} + D_{km} < (D_{im} + D_{jk} = D_{ik} + D_{jm})$$

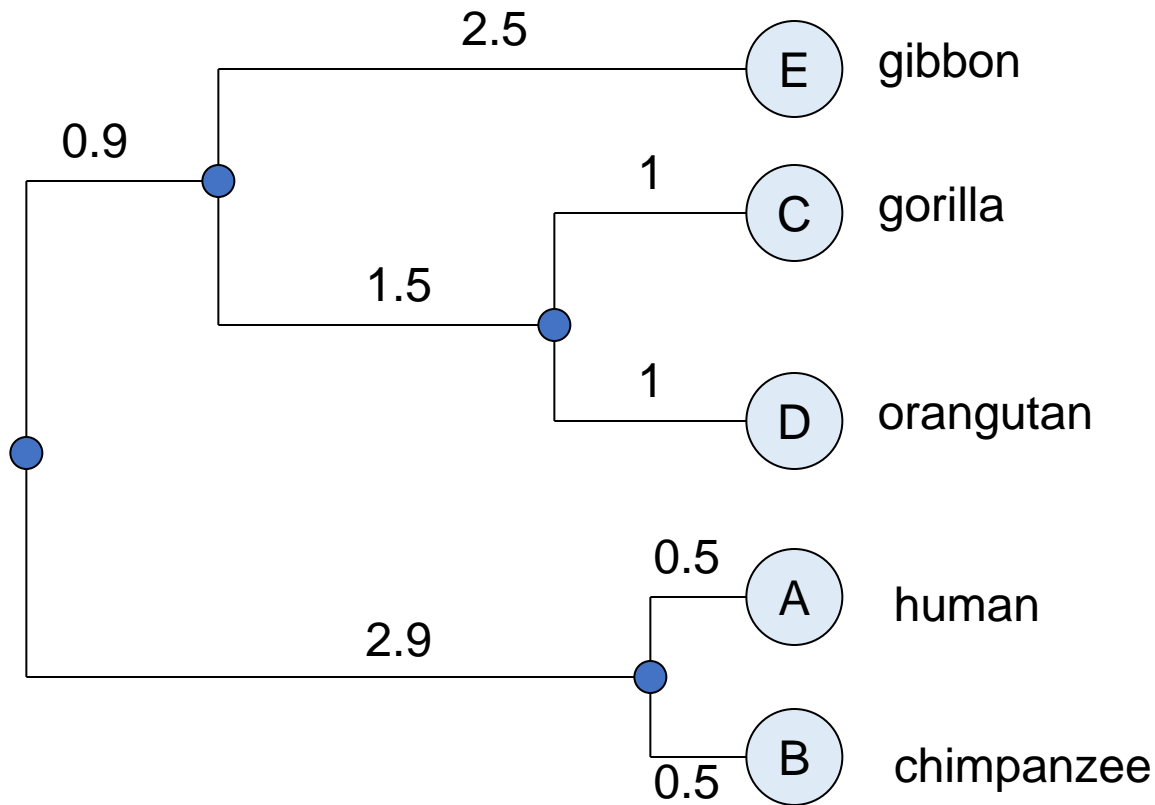
Without additive distances we **CANNOT** construct a phylogenetic tree

UPGMA algorithm - summary

- Initialization:
 - Create N clusters, 1 species per cluster
 - Set the size of each cluster to 1
 - Create leaf for each cluster
- Iterate (until only 1 cluster left)
 - Find C_i and C_j with min $d_{C_i C_j}$
 - Create a new cluster $C(ij)$ which has $n(ij)=n_i+n_j$ members
 - Connect C_i and C_j through a new parent node and set the distance from this new parent node to the leaf node of each cluster to $\frac{1}{2} d_{C_i C_j}$
 - Delete columns and rows that correspond to amalgamed clusters i and j
 - Add a column and a row for a new cluster
 - Compute distances from a new cluster $C(ij)$ to all remaining clusters:

$$d_{C(ij)C_k} = \sum_{\text{all } x \in C(ij), \text{ all } y \in C_k} d_{XY} / (n(ij) * n_k)$$

Ultrametric trees and UPGMA



This tree is clocklike, ultrametric: the total length of the path from a given internal node to each leaf is the same.

Assumption: the molecular clock of mutations ticks with a constant pace

- UPGMA reconstructs the tree based on the molecular clock assumption, that is why a new node is always created at the same distance from all the leaves

When the tree reflects reality

$$SSQ(T) = \sum_{i \text{ from } 1 \text{ to } N} \sum_{j \neq i} w_{ij} (d_{ij} - \text{Tree}D_{ij})$$

- If the solution to $SSQ(T)=0$, and there was a molecular clock with constant pace, then UPGMA guarantees to find the correct solution
- If not:
 - It can find a good enough solution, but **the correctness of the tree topology is not guaranteed**
 - Use *the neighbor-joining algorithm* to check the correctness of the tree topology. This algorithm relies on additivity but does not require the distances to be ultrametric

Test for ultrametric condition

- We can predict whether the reconstruction of the real tree is likely to be correct by testing our distances for *ultrametric condition*:

The distance matrix is *ultrametric* if for any triplet of sequences, X_i , X_j , X_k , the distances d_{ij} , d_{ik} , d_{jk} are either all equal or two are equal and the remaining one is smaller

Thus, if distances were derived from a real tree with a molecular clock, the distance matrix will be ultrametric