

CS 466

Introduction to Bioinformatics

Lecture 15

Mohammed El-Kebir

October 16, 2019



Course Announcements

HW 3 will be released Oct 25 – due Nov 1 by 11:59pm

Project proposal due on Nov 3
(Motivation, Datasets/papers, Planned method/experiments, Timeline)

Project report due on Dec 22

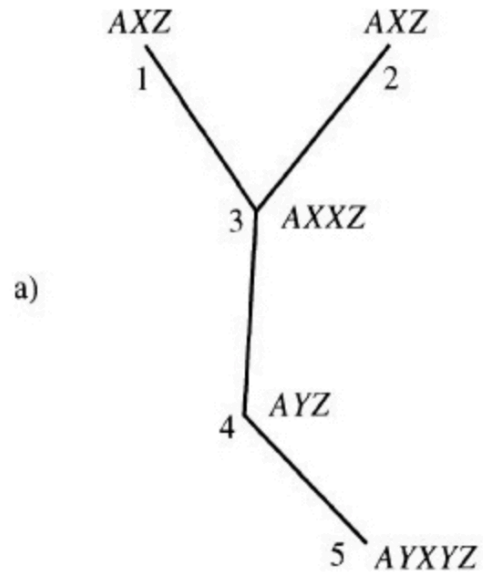
Outline

- Introduction
- Hierarchical clustering
- Additive distance phylogeny
- Four point condition
- Neighbor joining

Reading:

- Chapter 10.2 and 10.5-10.8 in Jones and Pevzner

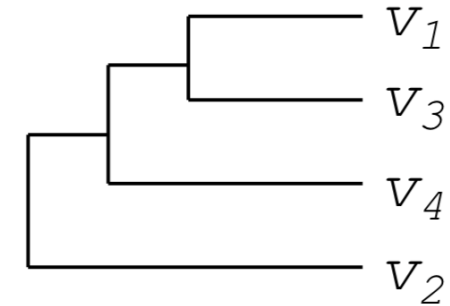
Alignments and Trees



Tree / star alignment

b)

3	A	X	X	_	Z
1	A	X	_	_	Z
2	A	_	X	_	Z
4	A	Y	_	_	Z
5	A	Y	X	X	Z



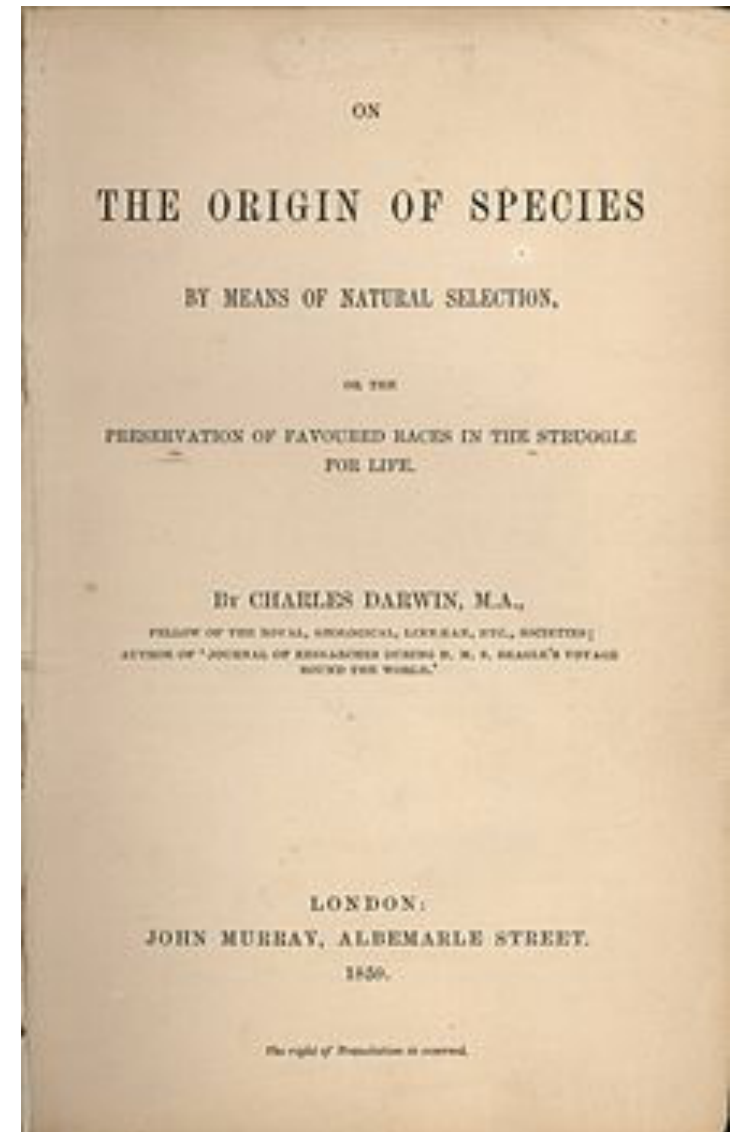
Guide tree in
progressive alignment

Tree topology represents similarity/distance between sequences

Biological sequences typically come from the present

Evolutionary Studies and Phylogenies

- Since Darwin's book (1859) until 1960s:
Phylogeny reconstruction from
anatomical features
- Subjective observations led to
inconclusive/incorrect phylogenies

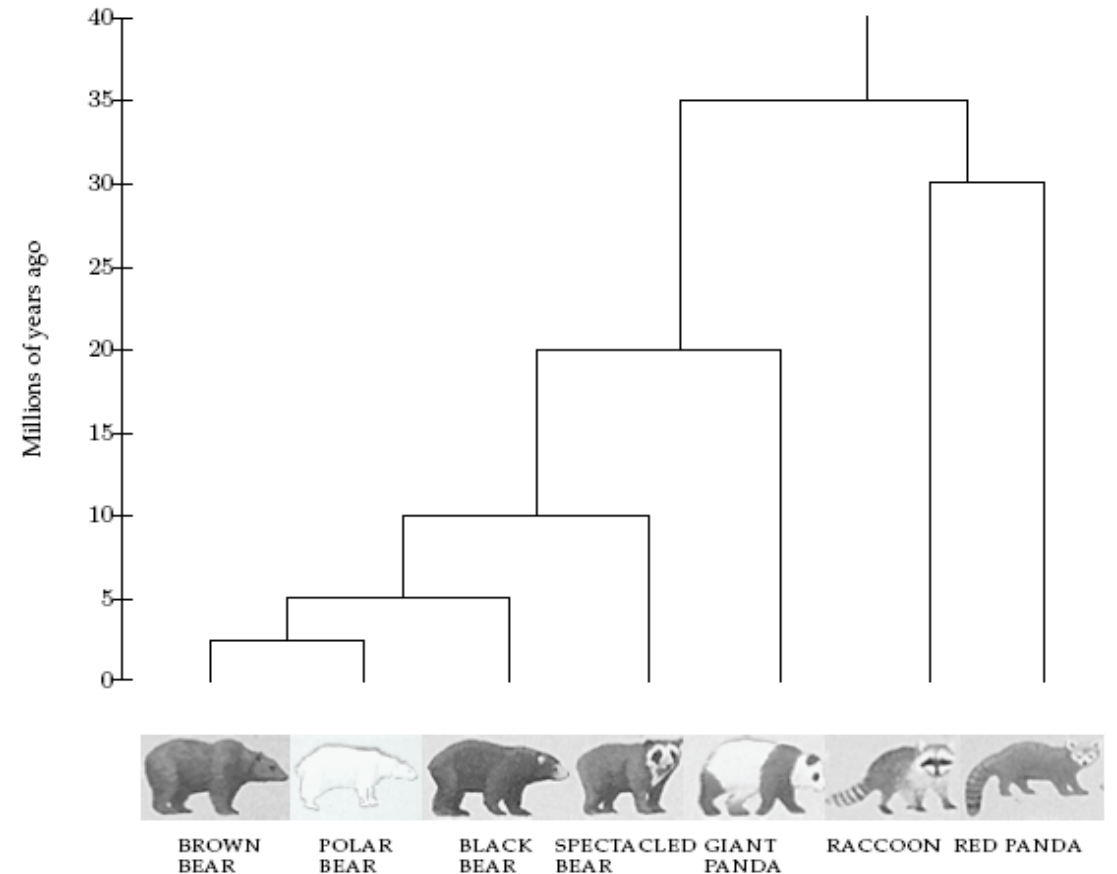


Evolutionary Studies and Phylogenies

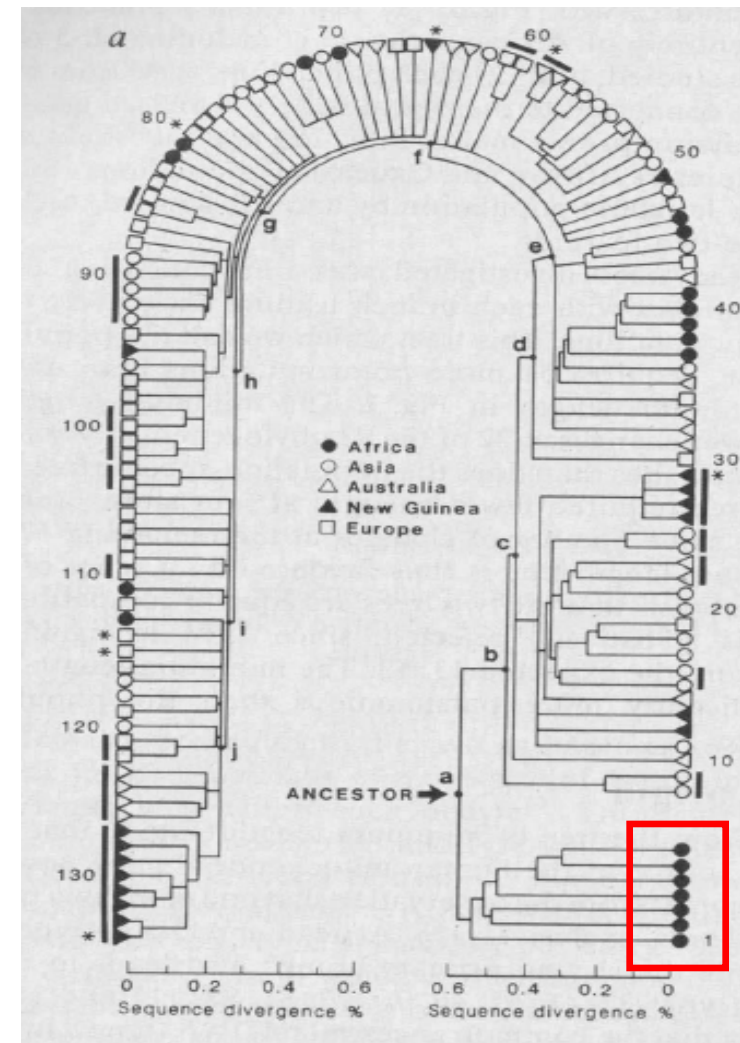
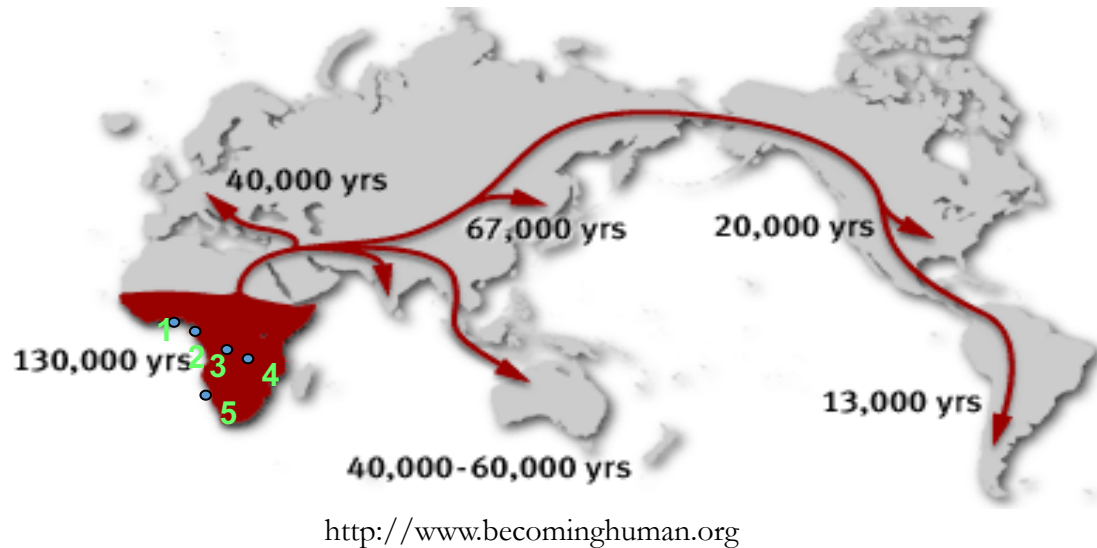
- Subjective observations led to inconclusive/incorrect phylogenies

Example

- Giant pandas look like bears but have features that are unusual for bears and typical for raccoons
- In 1985, Steven O'Brien and colleagues solved the giant panda classification problem using DNA sequences and algorithms



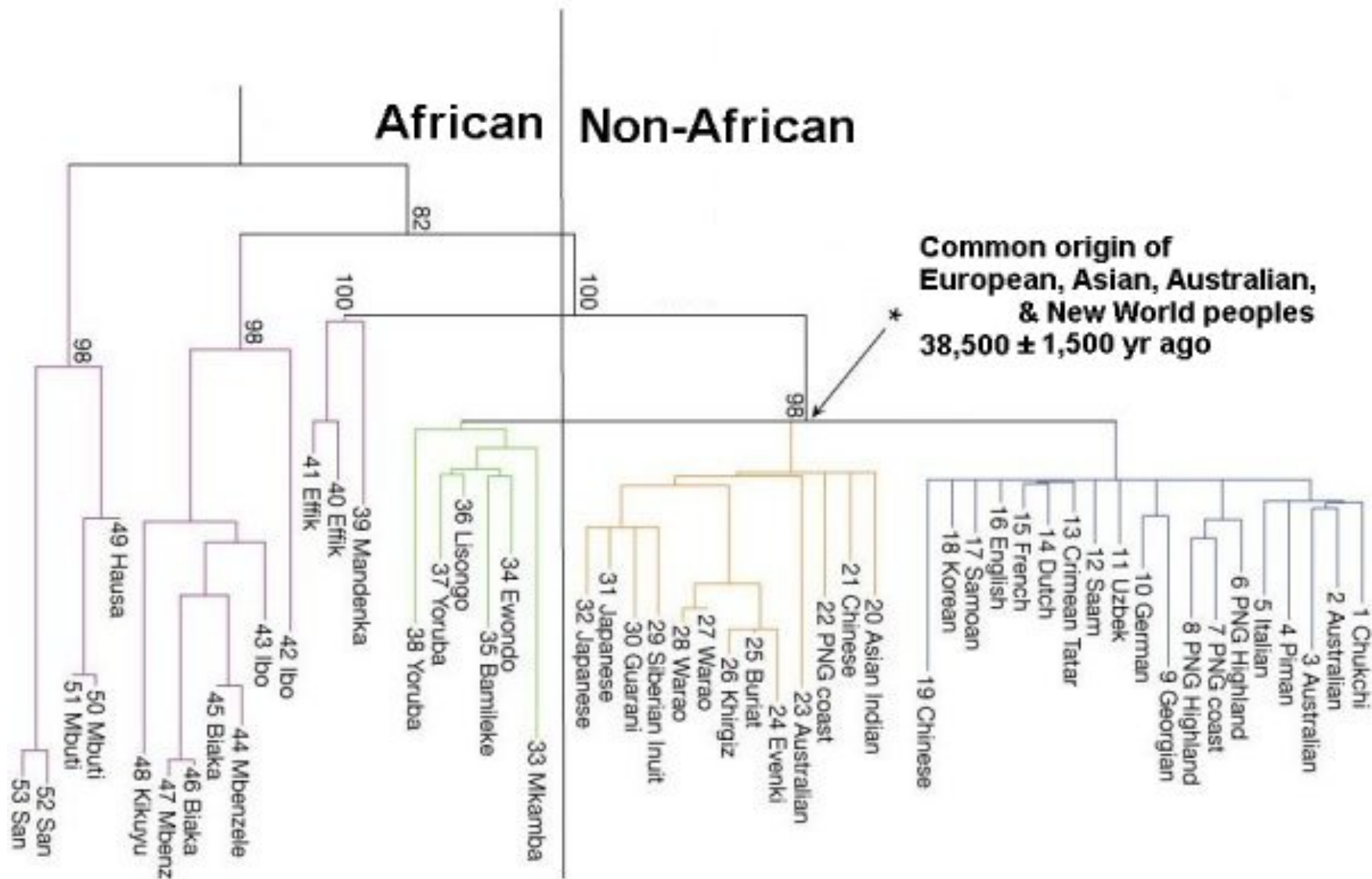
Out of Africa Hypothesis



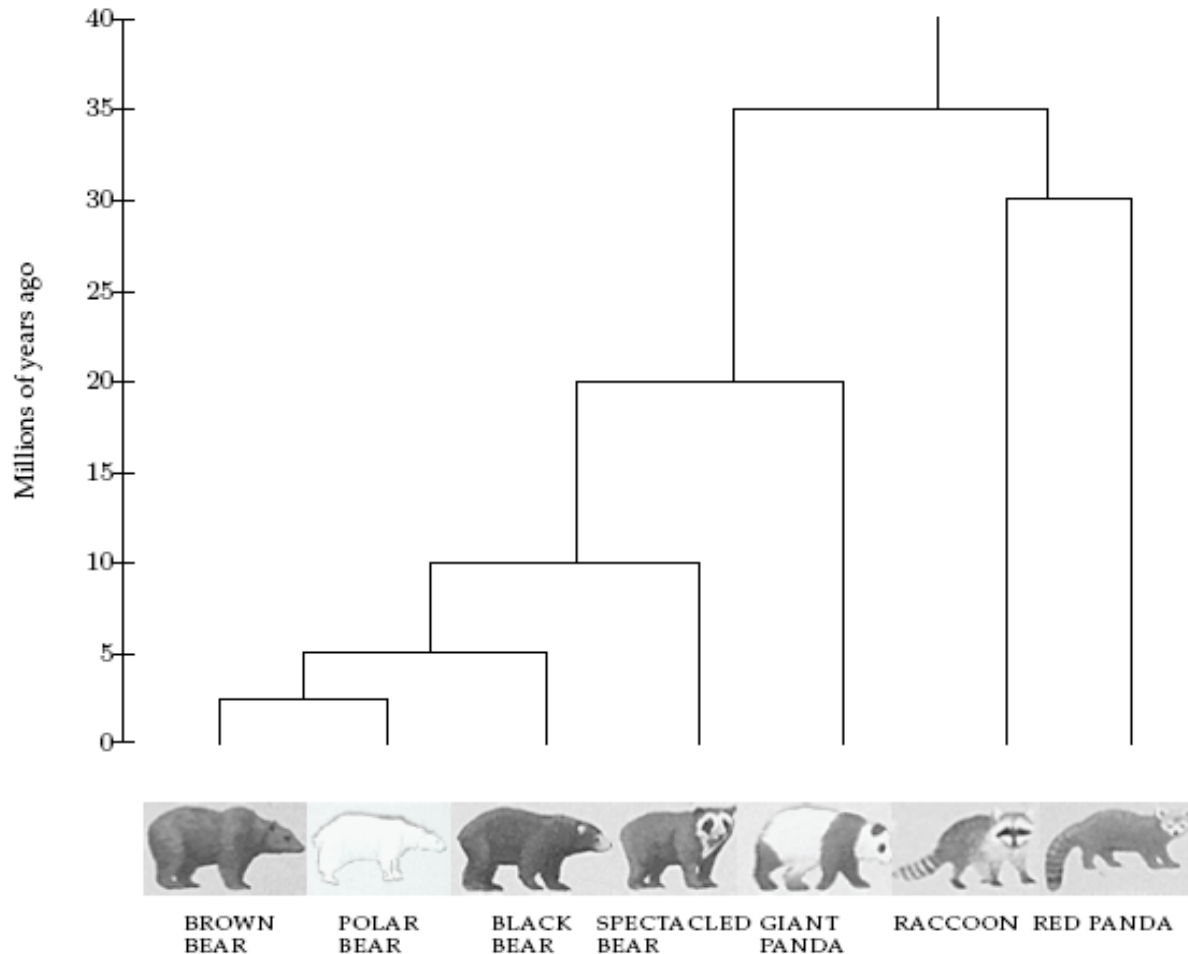
Vigilant, Stoneking, Harpending, Hawkes, and Wilson (1991)

Out of Africa Hypothesis claims that our most ancient ancestor lived in Africa roughly 200,000 years ago

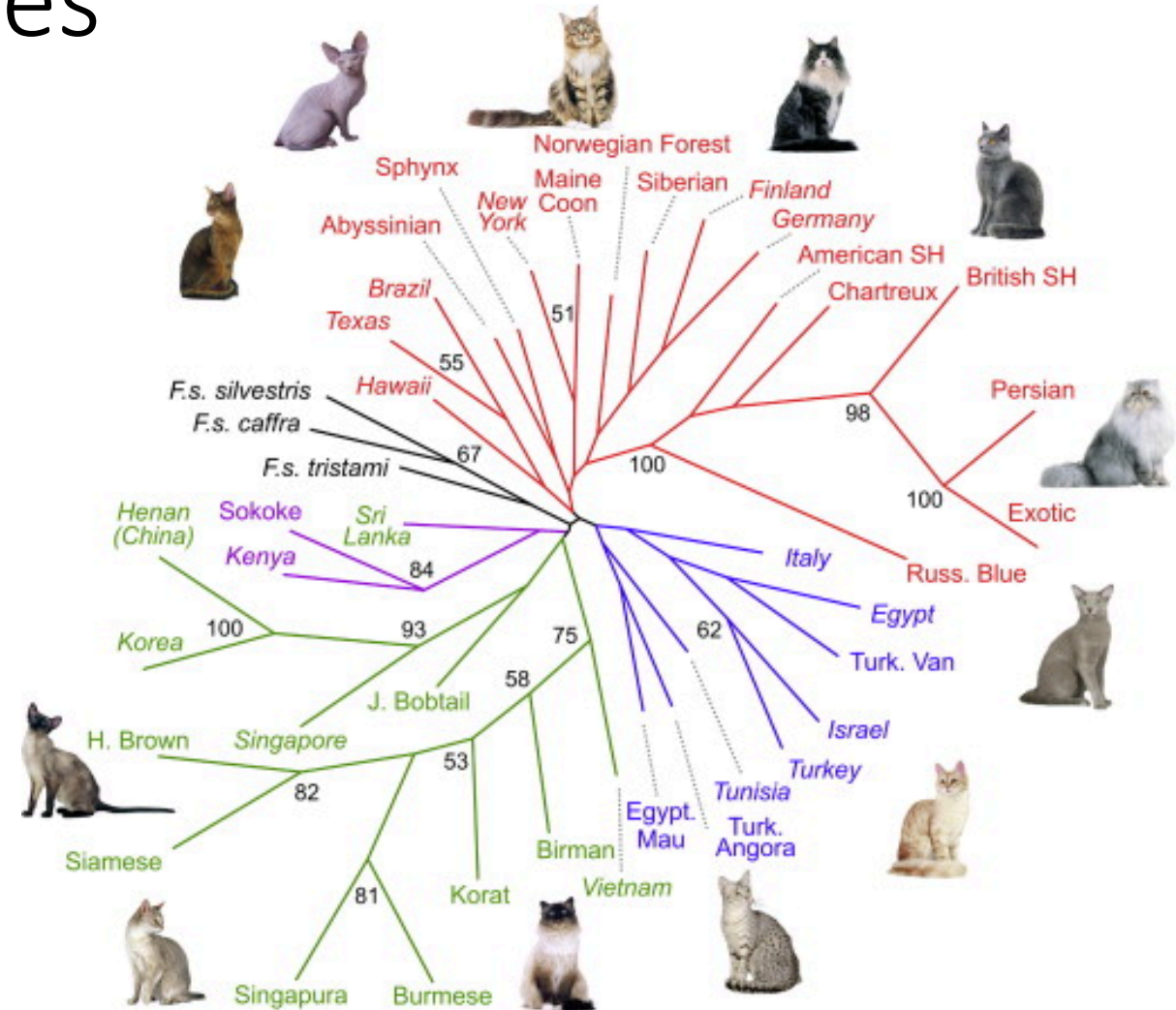
Evolutionary Tree of Humans



Evolutionary Tree of Species



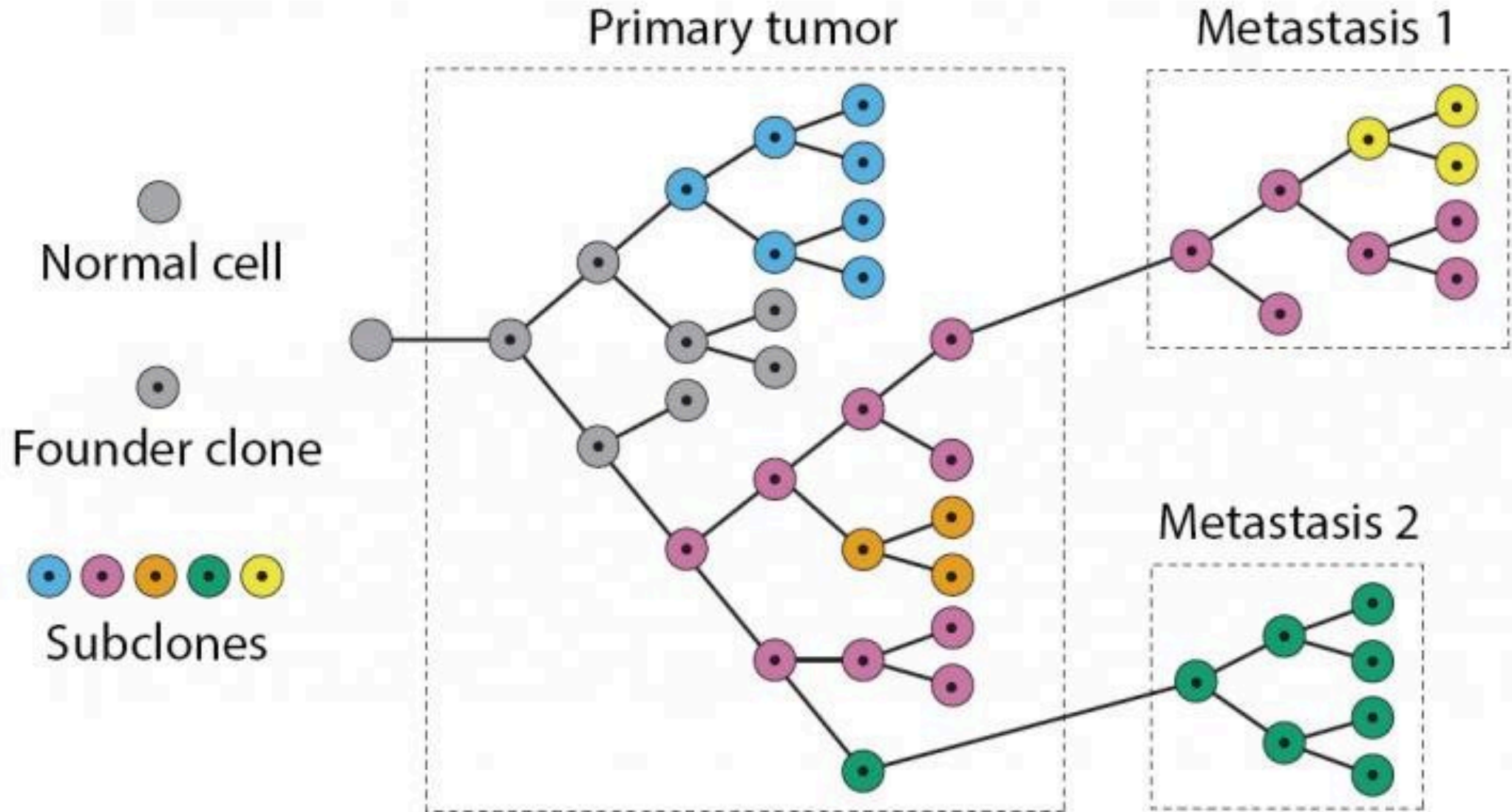
<http://bix.ucsd.edu/bioalgorithms/>



[Lipinski *et al.*, 2008]

Question: What are the evolutionary relationships between species?

Evolutionary Tree of a Tumor



<https://www.sciencedaily.com/releases/2016/09/160909223504.htm>

Phylogenetic Tree Reconstruction

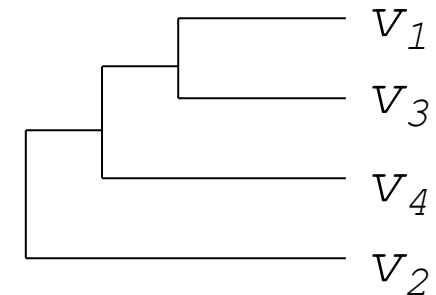
Mouse: ACAGTGACGCCACACACGT
Gorilla: CCTGTGACGTAACAAACGA
Chimpanzee: CCTGTGAGGTAGCAAACGA
Human: CCTGTGAGGTAGCACACGA

Distance Metric ↓

	V_1	V_2	V_3	V_4
V_1	—			
V_2	.17	—		
V_3	.87	.28	—	
V_4	.59	.33	.62	—

Distance Table

???



Phylogenetic Tree

Question: Given sequence data, how to reconstruct tree?

Outline

- Introduction
- Hierarchical clustering
- Additive distance phylogeny
- Four point condition
- Neighbor joining

Reading:

- Chapter 10.2 and 10.5-10.8 in Jones and Pevzner

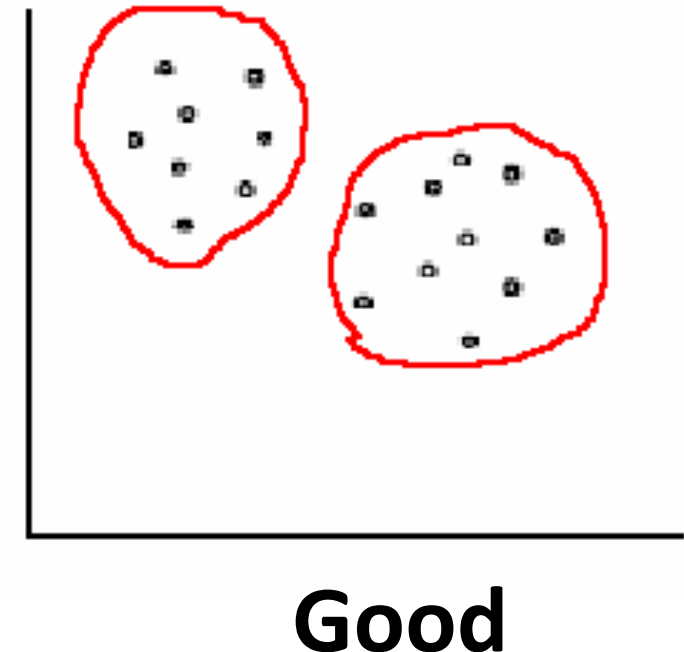
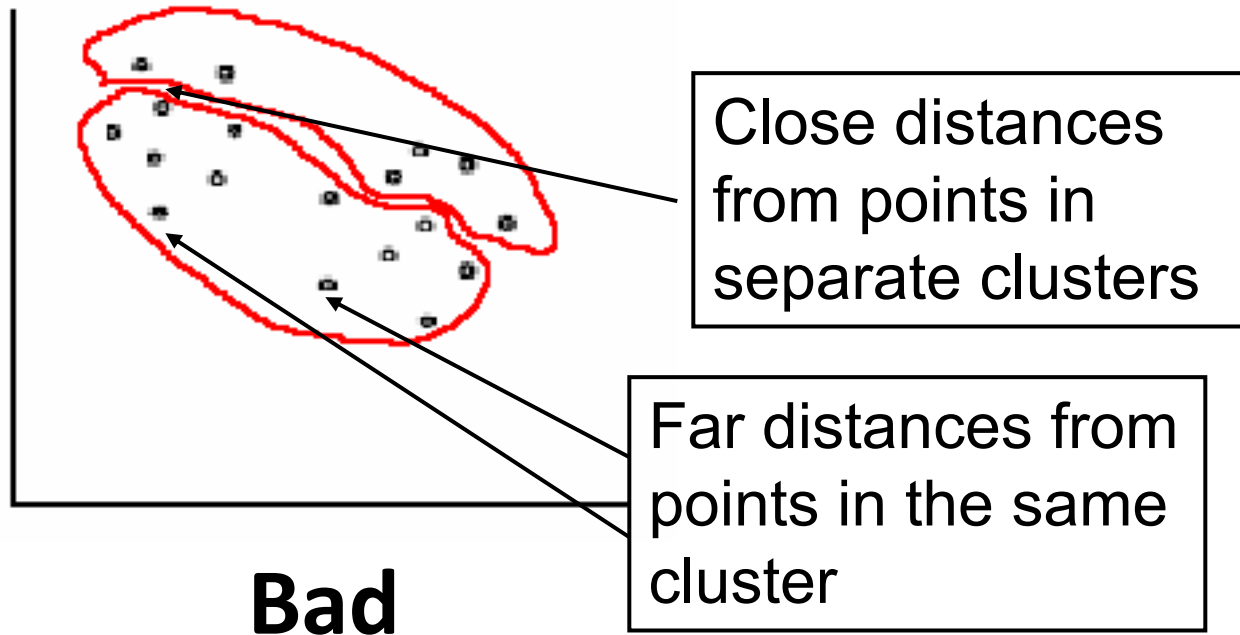
Clustering

Given:

(1) $n \times n$ matrix $D = [d_{i,j}]$

Want:

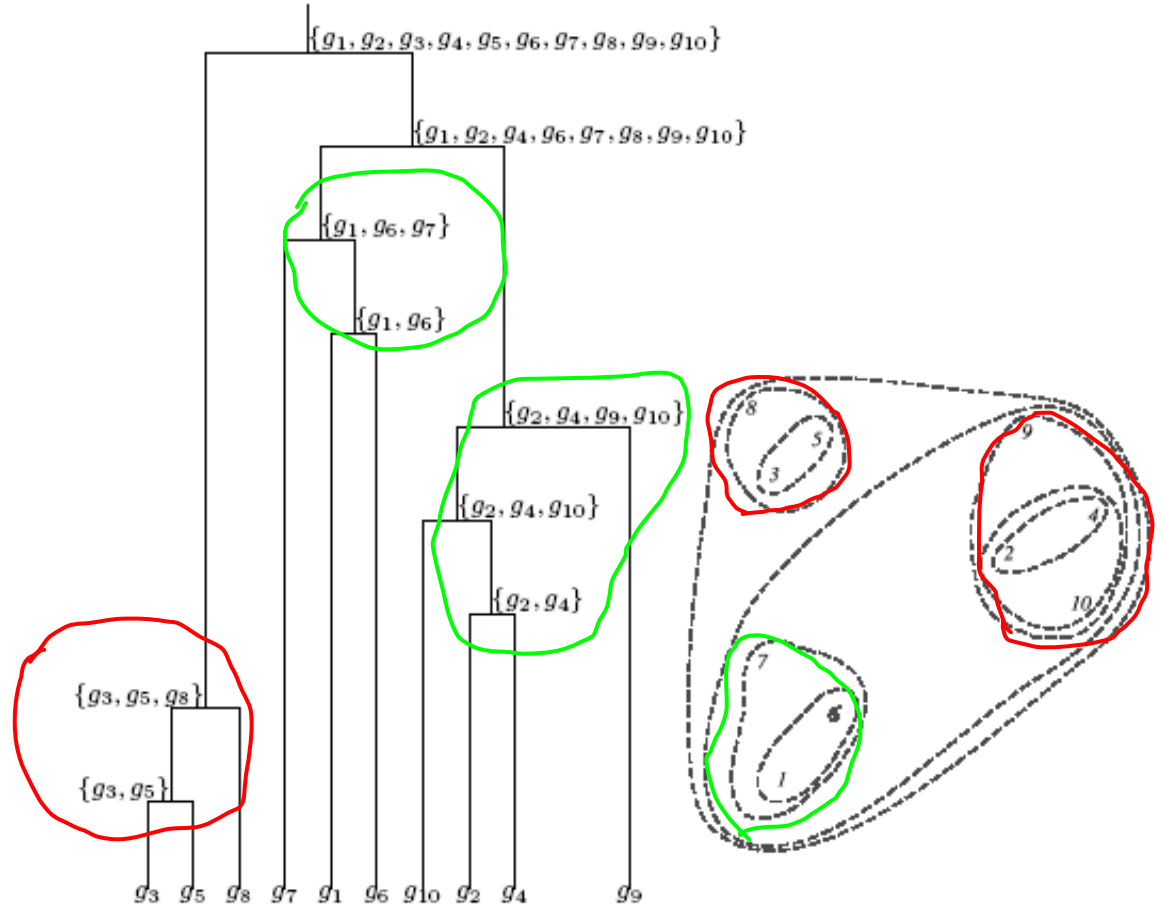
- (1) Homogeneity within clusters
- (2) Separation between clusters



Hierarchical Clustering

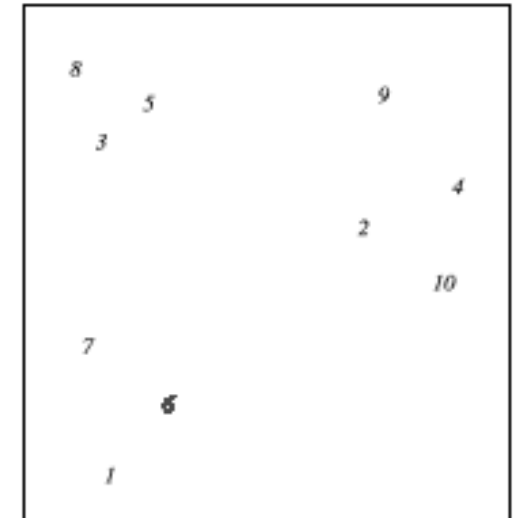
Organize elements into a tree such that:

- Leaves are elements
- Paths between leaves represent pairwise element distance
- Similar elements lie within same subtrees



Hierarchical Clustering

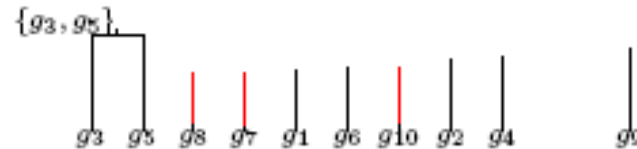
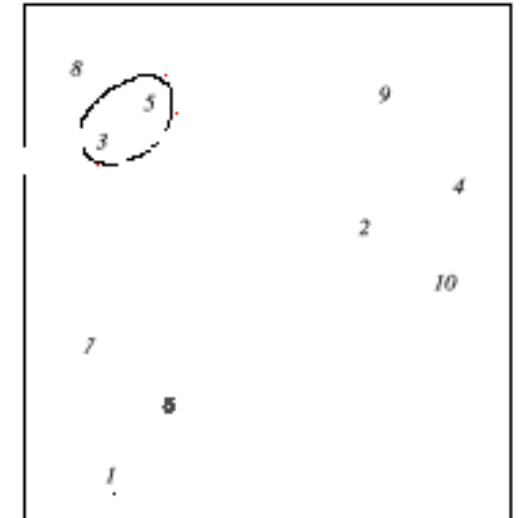
1. Hierarchical Clustering (\mathbf{D} , n)
2. Form n clusters each with one element
3. Construct a graph \mathbf{T} by assigning one vertex to each cluster
4. **while** there is more than one cluster
5. Find the two closest clusters C_1 and C_2
6. Merge C_1 and C_2 into new cluster C with $|C_1| + |C_2|$ elements
7. **Compute distance from C to all other clusters**
8. Add a new vertex \mathbf{C} to \mathbf{T} and connect to vertices C_1 and C_2
9. Remove rows and columns of \mathbf{D} corresponding to C_1 and C_2
10. Add a row and column to \mathbf{D} corresponding to the new cluster \mathbf{C}
11. **return \mathbf{T}**



g_3 g_5 g_8 g_7 g_1 g_6 g_{10} g_2 g_4 g_9

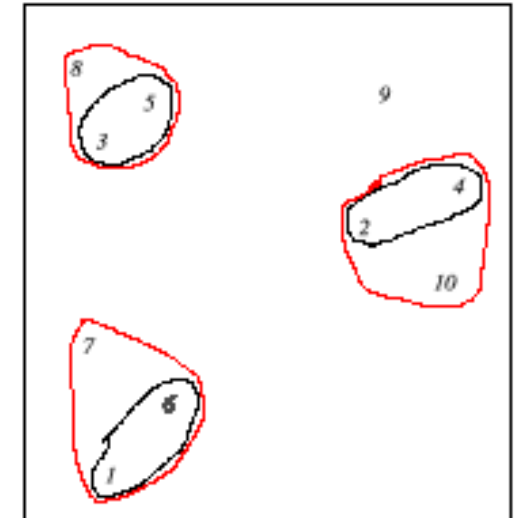
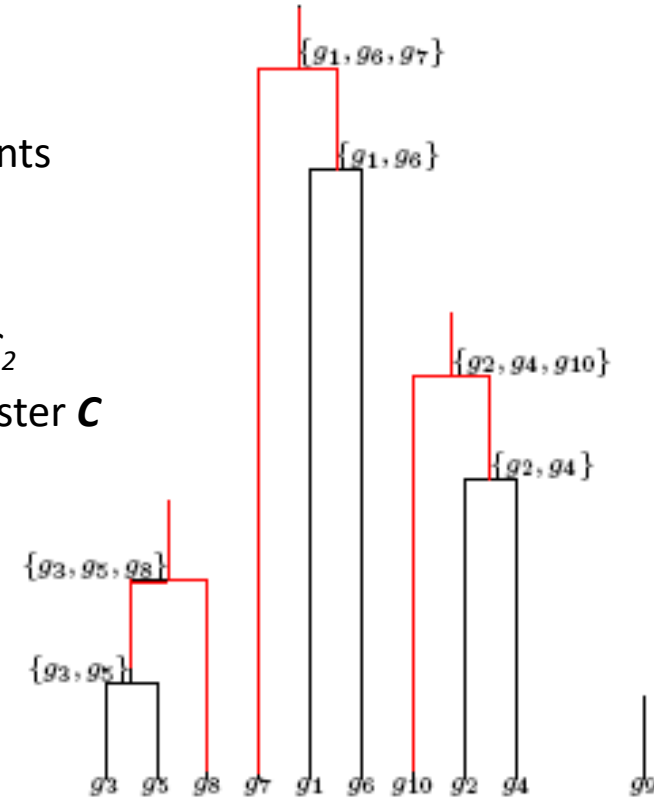
Hierarchical Clustering

1. Hierarchical Clustering (\mathbf{D} , n)
2. Form n clusters each with one element
3. Construct a graph \mathbf{T} by assigning one vertex to each cluster
4. **while** there is more than one cluster
5. Find the two closest clusters C_1 and C_2
6. Merge C_1 and C_2 into new cluster C with $|C_1| + |C_2|$ elements
7. **Compute distance from C to all other clusters**
8. Add a new vertex \mathbf{C} to \mathbf{T} and connect to vertices C_1 and C_2
9. Remove rows and columns of \mathbf{D} corresponding to C_1 and C_2
10. Add a row and column to \mathbf{D} corresponding to the new cluster \mathbf{C}
11. **return \mathbf{T}**



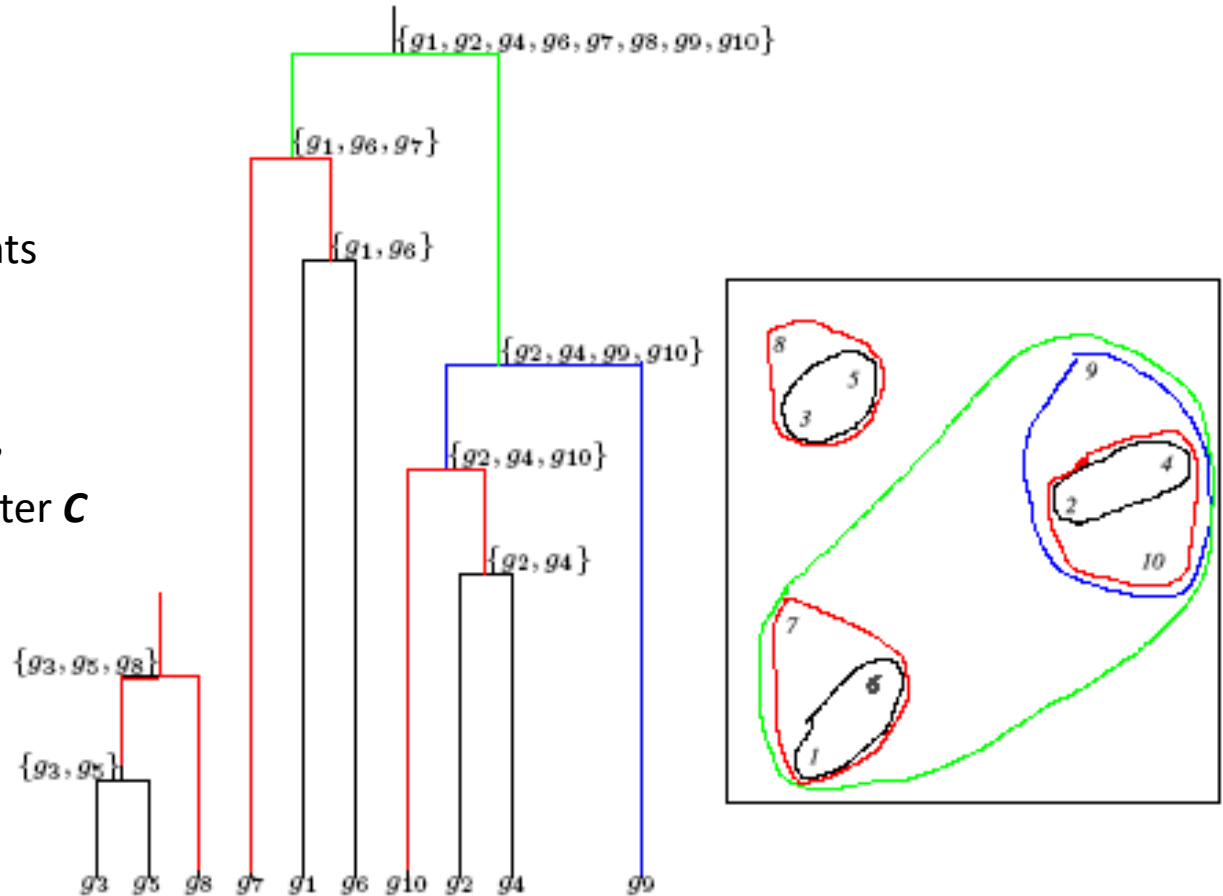
Hierarchical Clustering

1. Hierarchical Clustering (\mathbf{D} , n)
2. Form n clusters each with one element
3. Construct a graph \mathbf{T} by assigning one vertex to each cluster
4. **while** there is more than one cluster
5. Find the two closest clusters C_1 and C_2
6. Merge C_1 and C_2 into new cluster C with $|C_1| + |C_2|$ elements
7. **Compute distance from C to all other clusters**
8. Add a new vertex \mathbf{C} to \mathbf{T} and connect to vertices C_1 and C_2
9. Remove rows and columns of \mathbf{D} corresponding to C_1 and C_2
10. Add a row and column to \mathbf{D} corresponding to the new cluster \mathbf{C}
11. **return** \mathbf{T}



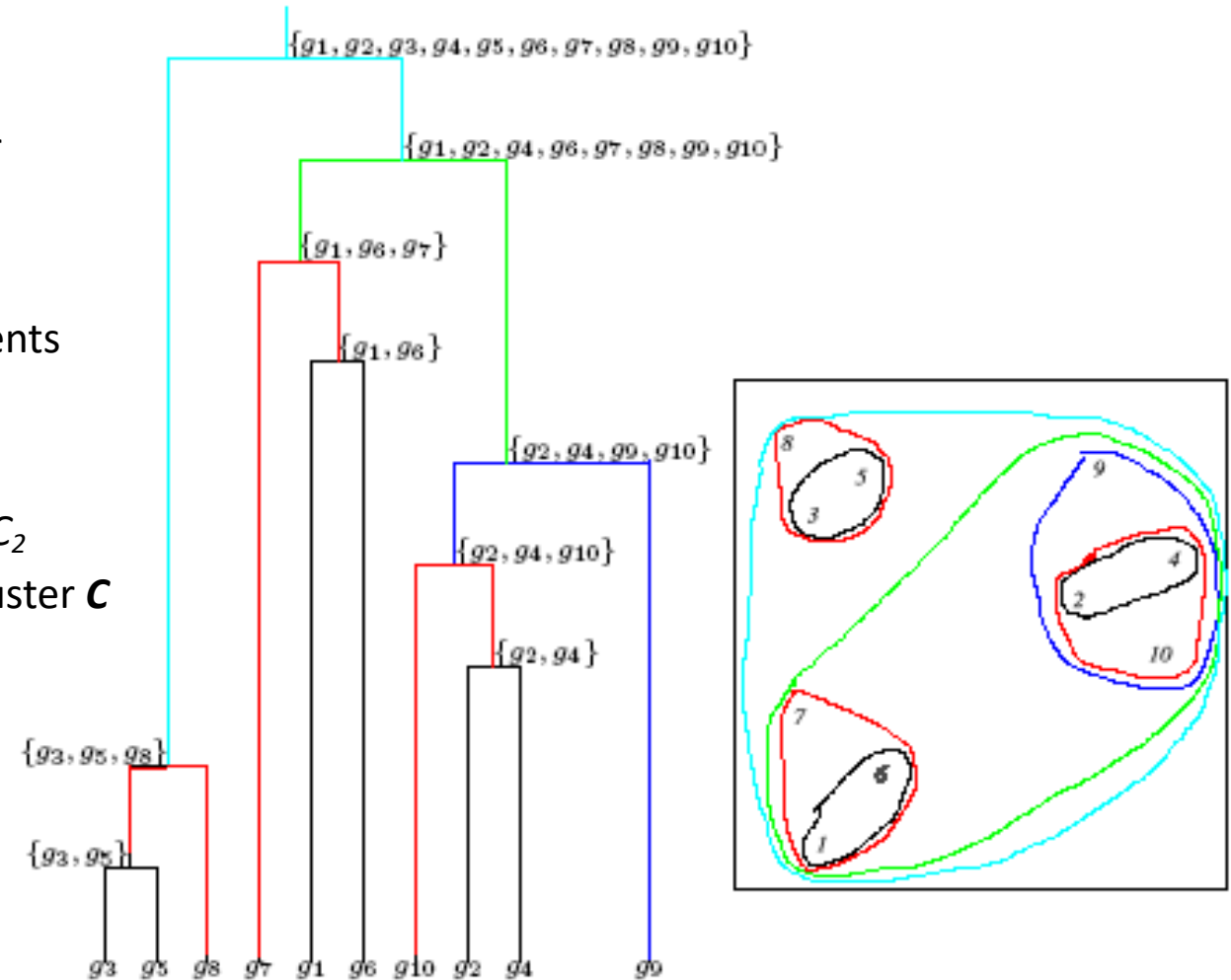
Hierarchical Clustering

1. Hierarchical Clustering (\mathbf{D}, n)
2. Form n clusters each with one element
3. Construct a graph \mathbf{T} by assigning one vertex to each cluster
4. **while** there is more than one cluster
5. Find the two closest clusters C_1 and C_2
6. Merge C_1 and C_2 into new cluster C with $|C_1| + |C_2|$ elements
7. **Compute distance from C to all other clusters**
8. Add a new vertex \mathbf{C} to \mathbf{T} and connect to vertices C_1 and C_2
9. Remove rows and columns of \mathbf{D} corresponding to C_1 and C_2
10. Add a row and column to \mathbf{D} corresponding to the new cluster \mathbf{C}
11. **return \mathbf{T}**



Hierarchical Clustering

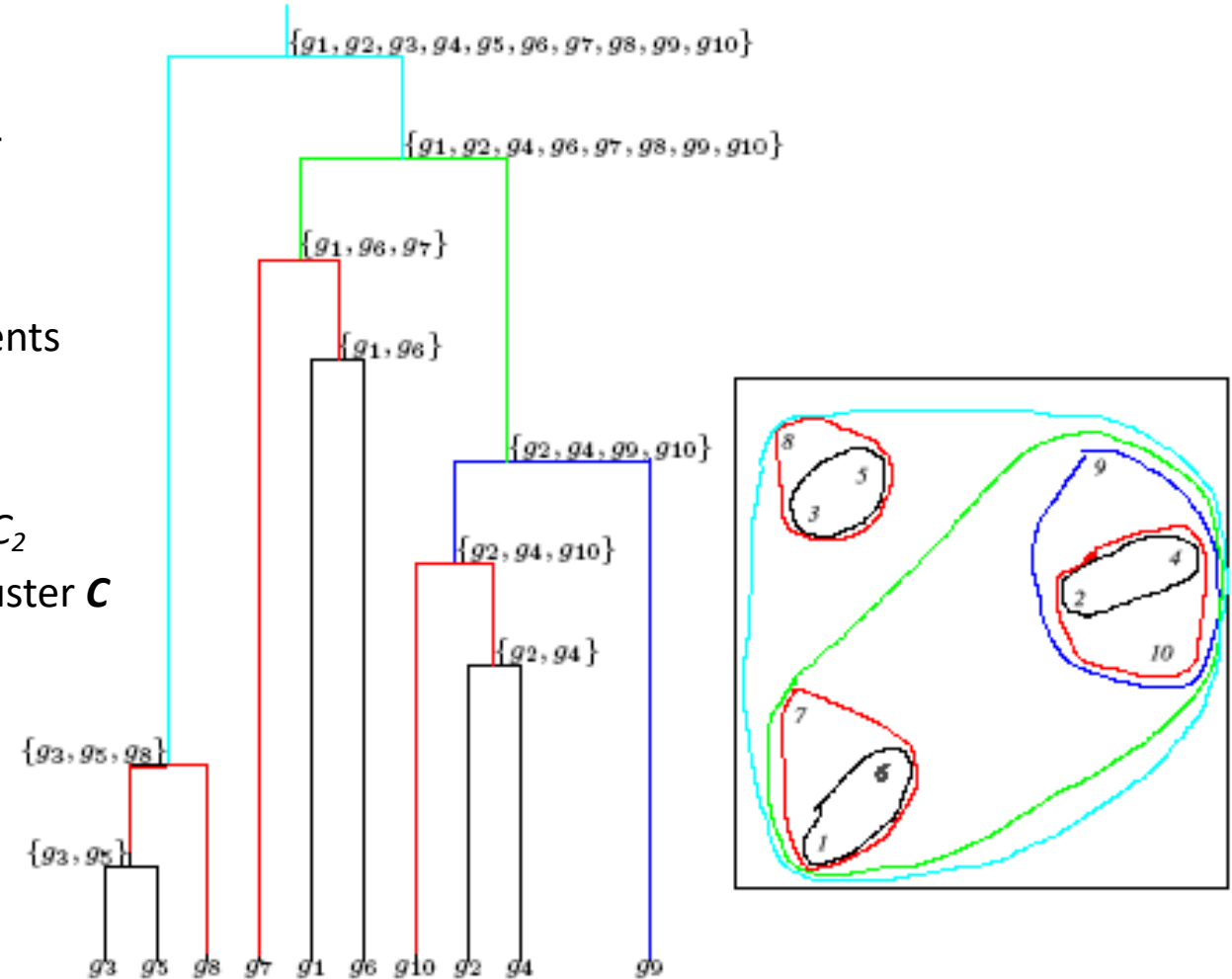
1. Hierarchical Clustering (\mathbf{D}, n)
2. Form n clusters each with one element
3. Construct a graph \mathbf{T} by assigning one vertex to each cluster
4. **while** there is more than one cluster
 5. Find the two closest clusters C_1 and C_2
 6. Merge C_1 and C_2 into new cluster C with $|C_1| + |C_2|$ elements
 7. **Compute distance from C to all other clusters**
 8. Add a new vertex \mathbf{C} to \mathbf{T} and connect to vertices C_1 and C_2
 9. Remove rows and columns of \mathbf{D} corresponding to C_1 and C_2
 10. Add a row and column to \mathbf{D} corresponding to the new cluster \mathbf{C}
11. return \mathbf{T}



Hierarchical Clustering

1. Hierarchical Clustering (\mathbf{D} , n)
2. Form n clusters each with one element
3. Construct a graph \mathbf{T} by assigning one vertex to each cluster
4. **while** there is more than one cluster
5. Find the two closest clusters C_1 and C_2
6. Merge C_1 and C_2 into new cluster C with $|C_1| + |C_2|$ elements
7. **Compute distance from C to all other clusters**
8. Add a new vertex \mathbf{C} to \mathbf{T} and connect to vertices C_1 and C_2
9. Remove rows and columns of \mathbf{D} corresponding to C_1 and C_2
10. Add a row and column to \mathbf{D} corresponding to the new cluster \mathbf{C}
11. **return \mathbf{T}**

Definition of distance
between clusters affects
clustering!



Hierarchical Clustering – Linkage Criteria

Names	Formula
Maximum or complete-linkage clustering	$\max \{ d(a, b) : a \in A, b \in B \}.$
Minimum or single-linkage clustering	$\min \{ d(a, b) : a \in A, b \in B \}.$
Mean or average linkage clustering, or UPGMA	$\frac{1}{ A \cdot B } \sum_{a \in A} \sum_{b \in B} d(a, b).$
Centroid linkage clustering, or UPGMC	$\ c_s - c_t\ $ where c_s and c_t are the centroids of clusters s and t , respectively.
Minimum energy clustering	$\frac{2}{nm} \sum_{i,j=1}^{n,m} \ a_i - b_j\ _2 - \frac{1}{n^2} \sum_{i,j=1}^n \ a_i - a_j\ _2 - \frac{1}{m^2} \sum_{i,j=1}^m \ b_i - b_j\ _2$

https://en.wikipedia.org/wiki/Hierarchical_clustering#Linkage_criteria

Outline

- Introduction
- Hierarchical clustering
- Additive distance phylogeny
- Four point condition
- Neighbor joining

Reading:

- Chapter 10.2 and 10.5-10.8 in Jones and Pevzner

Phylogenetic Tree Reconstruction

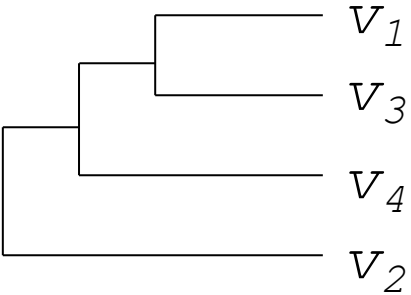
Mouse: ACAGTGACGCCACACACGT
Gorilla: CCTGTGACGTAACAAACGA
Chimpanzee: CCTGTGAGGTAGCAAACGA
Human: CCTGTGAGGTAGCACACGA

Distance Metric ↓

	V_1	V_2	V_3	V_4
V_1	—			
V_2	.17	—		
V_3	.87	.28	—	
V_4	.59	.33	.62	—

Distance Table

???



Phylogenetic Tree

Question: Given sequence data, how to reconstruct tree?

Distance

A **distance** (metric) on a set X is a function $d : X \times X \rightarrow \mathbb{R}$ s.t. for all $x, y, z \in X$:

- $i. \quad d(x, y) \geq 0$ [non-negativity]
- $ii. \quad d(x, y) = 0$ if and only if $x = y$ [identity of indiscernibles]
- $iii. \quad d(x, y) = d(y, x)$ [symmetry]
- $iv. \quad d(x, y) \leq d(x, z) + d(z, y)$ [triangle inequality]

Examples:

- $X = \mathbb{R}$ and $d(x, y) = |x - y|$
- $X = \Sigma^*$ and d is Hamming distance
- $X = \Sigma^*$ and d is edit distance

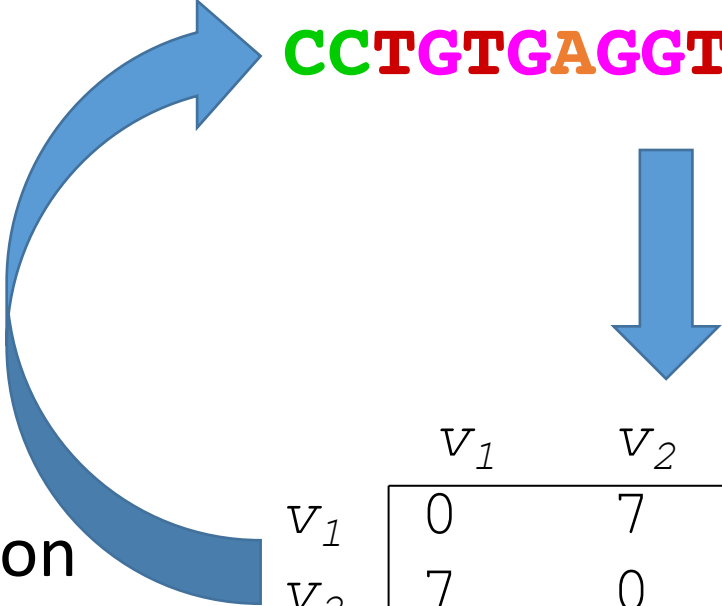
Alignment vs. Distance Matrices

Mouse: **ACAGTGACGCCACACACGT**
Gorilla: **CCTGTGACGTAACAAACGA**
Chimpanzee: **CCTGTGAGGTAGCAAACGA**
Human: **CCTGTGAGGTAGCACACGA**

Genes of length m in
 n species

Easy: use (weighted) edit distance

Reverse
transformation
not possible
due to loss of
information

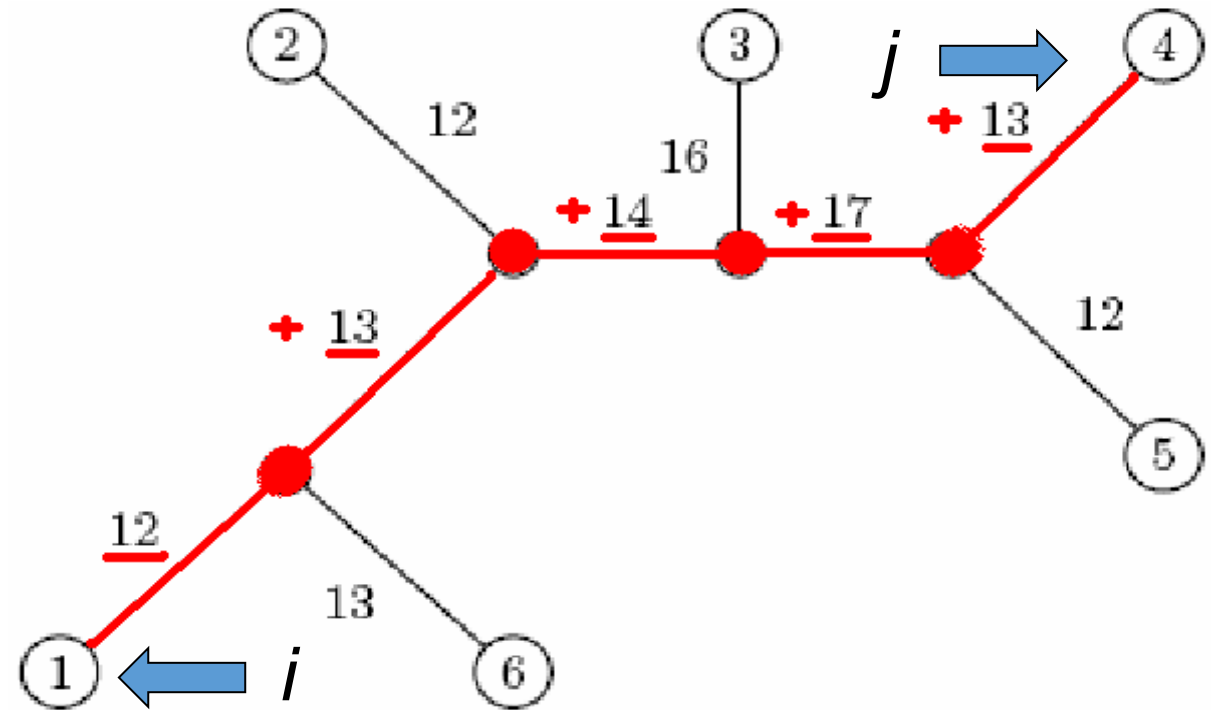


	v_1	v_2	v_3	v_4
v_1	0	7	11	10
v_2	7	0	4	6
v_3	11	4	0	2
v_4	10	6	2	0

$n \times n$ distance matrix

Distances in Trees

Given a tree T with positive edge weights $w(e)$, **tree distance** $d_T(i, j)$ between two leaves i and j is the sum of weights of edges on the unique path from i to j



$$d_T(1,4) = 12 + 13 + 14 + 17 + 13 = 69$$

General Distance vs. Tree Distance

Rat:

Mouse:

Gorilla:

Chimpanzee:

Human:

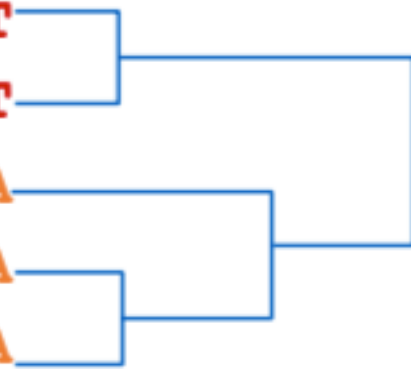
ACAGTGACGCCCCAAACGT

ACAGTGACGCTACAAACGT

CCTGTGACGTAACAAACGA

CCTGTGACGTAGCAAACGA

CCTGTGACGTAGCAAACGA



Tree distance $d_T(i, j)$ not necessarily equal to $d_{i, j}$ as given by distance matrix obtained from alignment

Fitting a Tree to a Given Distance Matrix

- Given n species, we can compute $n \times n$ distance matrix $D = [d_{i,j}]$
- Evolution of these n species is described by an unknown tree
- We need an algorithm to construct tree T that best fits D


Fitting a Tree to a Given Distance Matrix

- Given n species, we can compute $n \times n$ distance matrix $D = [d_{i,j}]$
- Evolution of these n species is described by an unknown tree
- We need an algorithm to construct tree T that best fits D

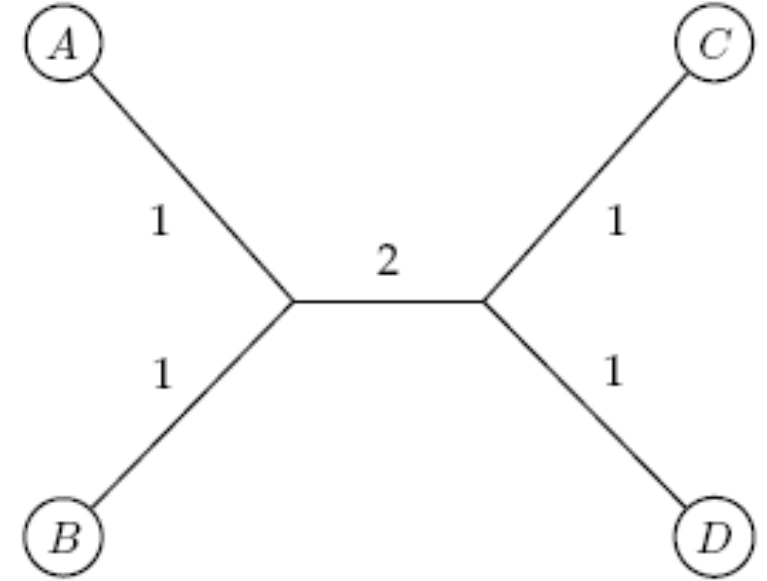
Distance-Based Phylogeny: Given $n \times n$ distance matrix $D = [d_{i,j}]$, find edge-weighted tree T with n leaves that best fits D

Question: How to define 'best fit'?


Additive Distance Matrices

Matrix D is  ADDITIVE if there exists a tree T with $d_{ij}(T) = D_{ij}$

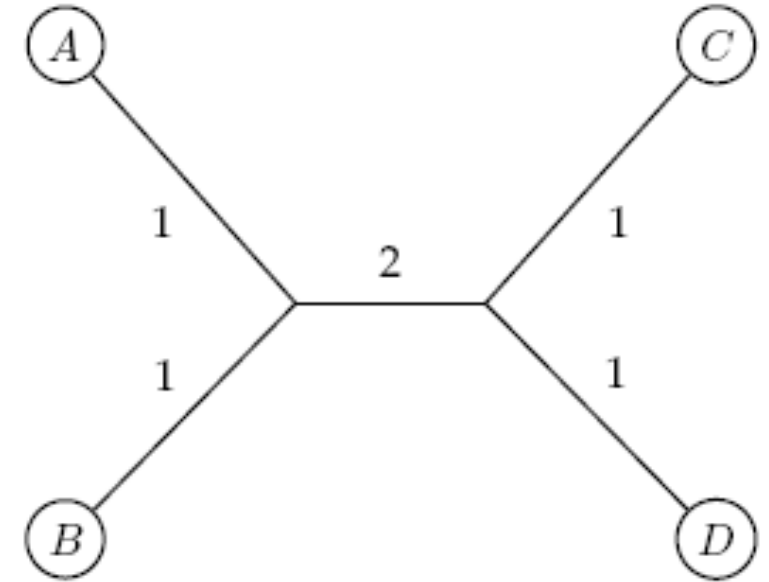
	A	B	C	D
A	0	2	4	4
B	2	0	4	4
C	4	4	0	2
D	4	4	2	0



Additive Distance Matrices

Matrix D is  ADDITIVE if there exists a tree T with $d_{ij}(T) = D_{ij}$

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



NON-ADDITIVE
otherwise 

	A	B	C	D
A	0	2	2	2
B	2	0	3	2
C	2	3	0	2
D	2	2	2	0

?

This is a constructive definition

A Small and a Large Problem

Small Additive Distance Phylogeny Problem:

Given $n \times n$ distance matrix $D = [d_{i,j}]$ and unweighted tree T with n leaves, determine edge weights such that $d_T(i, j) = d_{i,j}$

A Small and a Large Problem

Small Additive Distance Phylogeny Problem:

Given $n \times n$ distance matrix $D = [d_{i,j}]$ and unweighted tree T with n leaves, determine edge weights such that $d_T(i, j) = d_{i,j}$

Large Additive Distance Phylogeny Problem:

Given $n \times n$ distance matrix $D = [d_{i,j}]$, find tree T with n leaves **and** edge weights such that $d_T(i, j) = d_{i,j}$

A Small and a Large Problem

Small Additive Distance Phylogeny Problem:

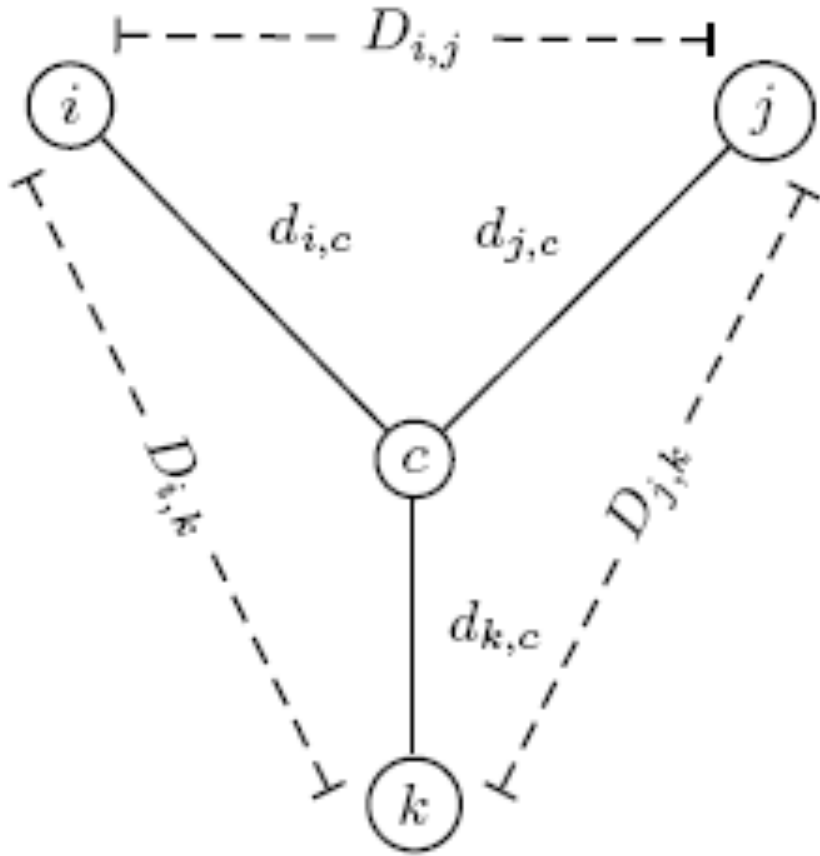
Given $n \times n$ distance matrix $D = [d_{i,j}]$ and unweighted tree T with n leaves, determine edge weights such that $d_T(i, j) = d_{i,j}$

Large Additive Distance Phylogeny Problem:

Given $n \times n$ distance matrix $D = [d_{i,j}]$, find tree T with n leaves **and** edge weights such that $d_T(i, j) = d_{i,j}$

Both problems can be solved in polynomial time

Additive Distance Problem with $n = 3$ Sequences



Additive Distance Problem with $n > 3$ Sequences

Unrooted binary tree with n leaves has $2n - 3$ edges and $\binom{n}{2}$ pairwise distances:

- $2n - 3$ variables
- $\binom{n}{2}$ equations

NON-ADDITIVE
otherwise 

	A	B	C	D
A	0	2	2	2
B	2	0	3	2
C	2	3	0	2
D	2	2	2	0

?

Solution not always possible for $n > 3$

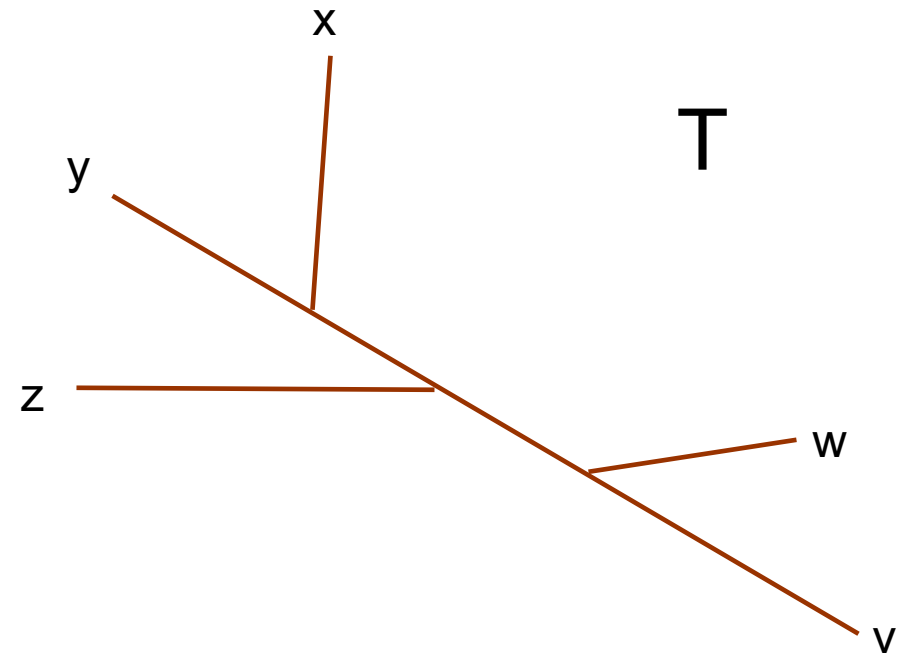
Small Additive Distance Problem

Small Additive Distance Phylogeny Problem:

Given $n \times n$ distance matrix $D = [d_{i,j}]$ and unweighted tree T with n leaves, determine edge weights such that $d_T(i, j) = d_{i,j}$

D

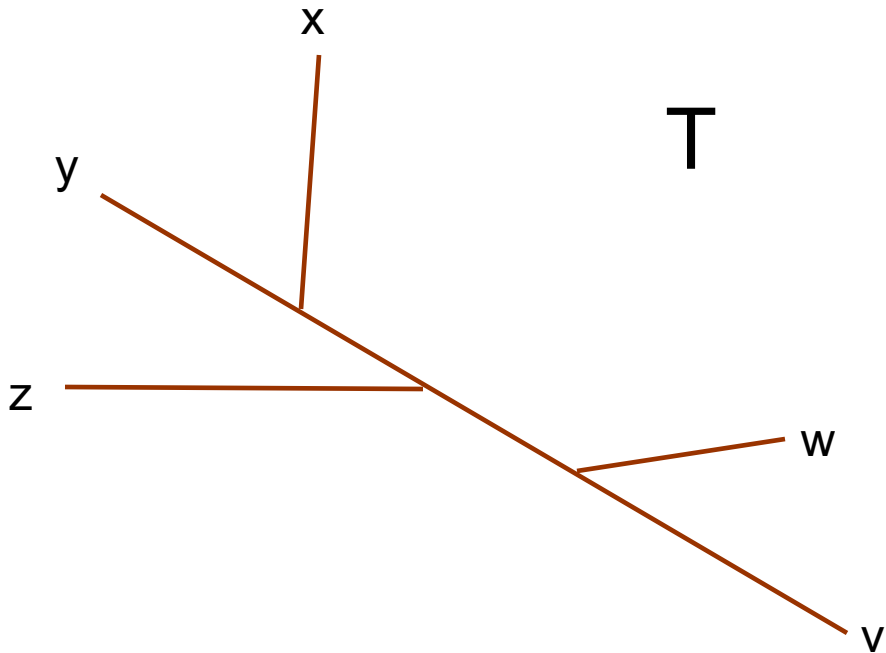
	v	w	x	y	z
v	0	10	17	16	16
w		0	15	14	14
x			0	9	15
y				0	14
z					0



Small Additive Distance Problem

D

	v	w	x	y	z
v	0	10	17	16	16
w		0	15	14	14
x			0	9	15
y				0	14
z					0



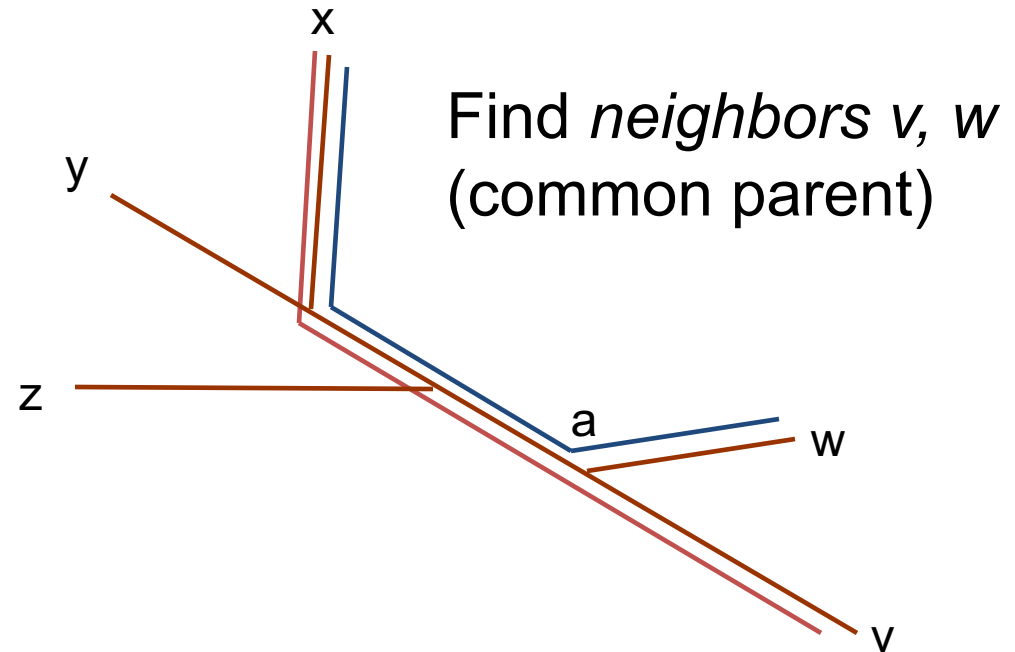
Small Additive Distance Problem

D

	v	w	x	y	z
v	0	10	17	16	16
w		0	15	14	14
x			0	9	15
y				0	14
z					0

D₁

	a	x	y	z
a	0	11	10	10
x		0	9	15
y			0	14
z				0



$$d_{ax} = \frac{1}{2} (d_{vx} + d_{wx} - d_{vw})$$

$$d_{ay} = \frac{1}{2} (d_{vy} + d_{wy} - d_{vw})$$

$$d_{az} = \frac{1}{2} (d_{vz} + d_{wz} - d_{vw})$$

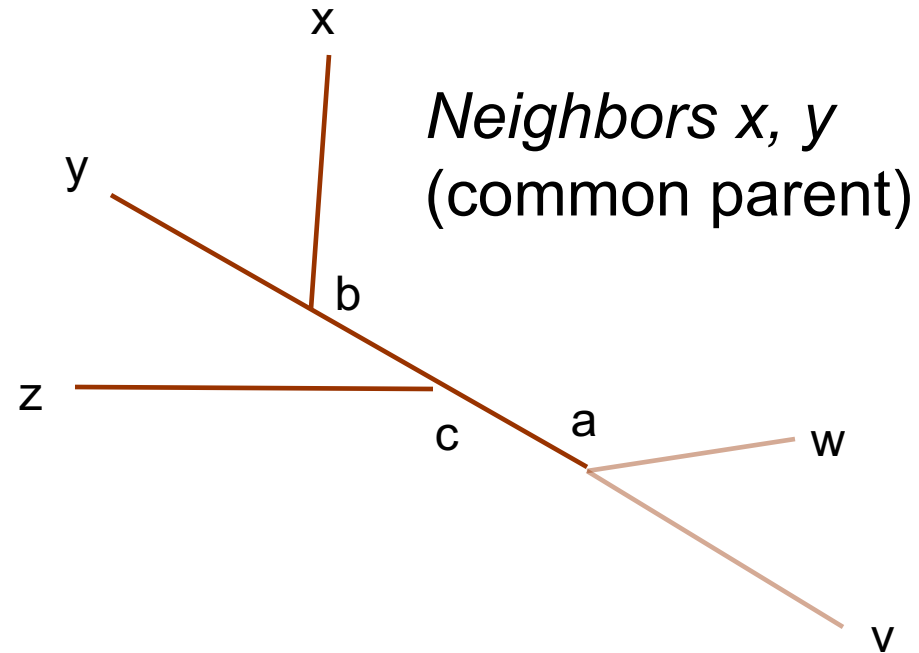
Small Additive Distance Problem

D_1

	a	x	y	z
a	0	11	10	10
x		0	9	15
y			0	14
z				0

D_2

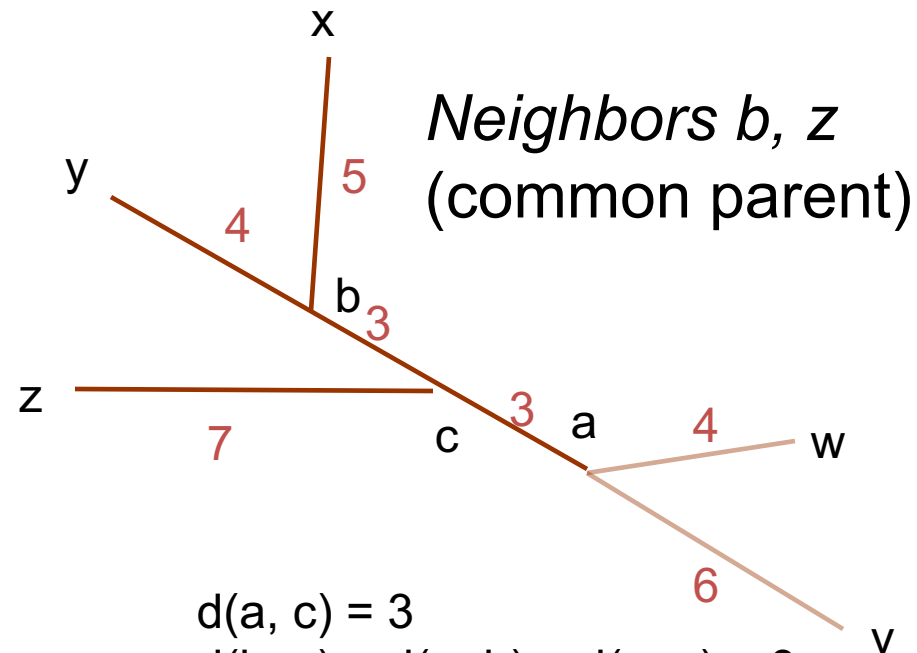
	a	b	z
a	0	6	10
b		0	10
z			0



Small Additive Distance Problem

D_1

	a	x	y	z
a	0	11	10	10
x		0	9	15
y			0	14
z				0



D_2

	a	b	z
a	0	6	10
b		0	10
z			0

D_3

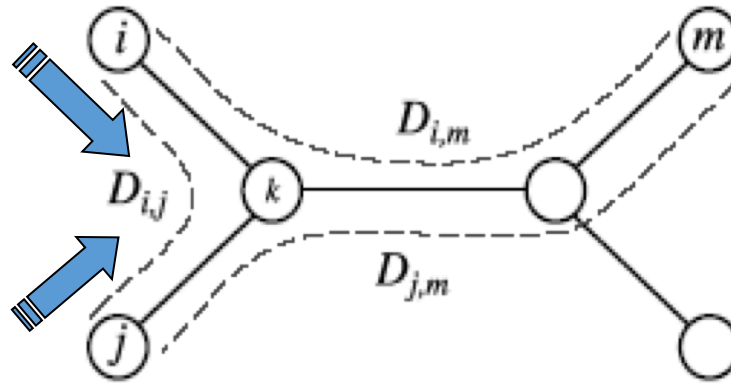
	a	c
a	0	3
c		0

$$\begin{aligned}
 d(a, c) &= 3 \\
 d(b, c) &= d(a, b) - d(a, c) = 3 \\
 d(c, z) &= d(a, z) - d(a, c) = 7 \\
 d(b, x) &= d(a, x) - d(a, b) = 5 \\
 d(b, y) &= d(a, y) - d(a, b) = 4 \\
 d(a, w) &= d(z, w) - d(a, z) = 4 \\
 d(a, v) &= d(z, v) - d(a, z) = 6
 \end{aligned}$$

Correct!!!

Small Additive Distance Problem

1. Find neighboring leaves i and j with parent k
2. Remove the rows and columns of i and j
3. Add a new row and column corresponding to k , where the distance from k to any other leaf m is computed as



$$d_{k,m} = \frac{(d_{i,m} + d_{j,m} - d_{i,j})}{2}$$

4. Repeat steps 1-3 until tree has only two vertices

A Small and a Large Problem

Small Additive Distance Phylogeny Problem:

Given $n \times n$ distance matrix $D = [d_{i,j}]$ and unweighted tree T with n leaves, determine edge weights such that $d_T(i, j) = d_{i,j}$

Large Additive Distance Phylogeny Problem:

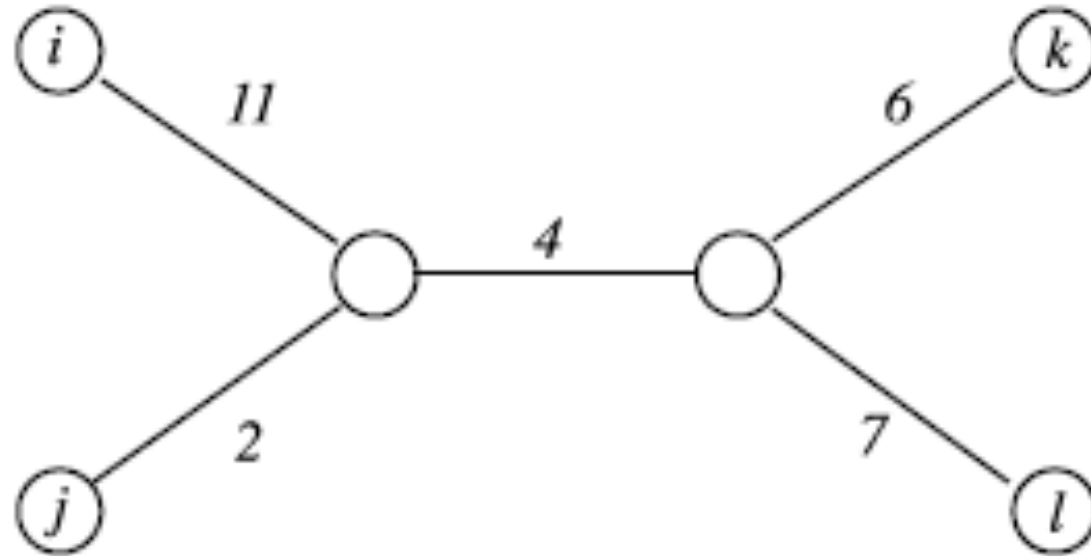
Given $n \times n$ distance matrix $D = [d_{i,j}]$, find tree T with n leaves **and** edge weights such that $d_T(i, j) = d_{i,j}$

Both problems can be solved in polynomial time

Large Additive Distance Phylogeny Problem

Idea: find neighboring leaves by simply selecting pair of closest leaves

	<i>i</i>	<i>j</i>	<i>k</i>	<i>l</i>
<i>i</i>	0	13	21	22
<i>j</i>		0	12	13
<i>k</i>			0	13
<i>l</i>				0



WRONG!

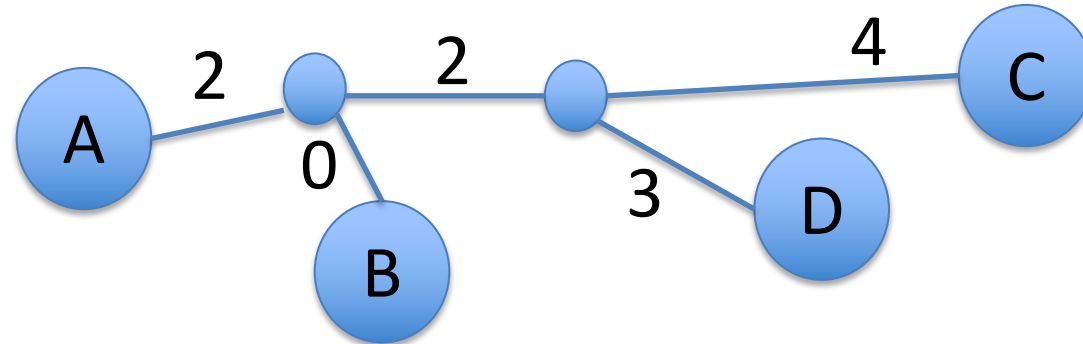
i and *j* are neighbors, but $(d_{ij} = 13) > (d_{jk} = 12)$.

Finding a pair of neighboring leaves is a nontrivial problem!

Degenerate Triples

A **degenerate triple** is a set of three distinct elements $i, j, k \in [n]$ such that $d_{i,j} + d_{j,k} = d_{i,k}$

	A	B	C	D
A	0	2	8	7
B		0	6	5
C			0	7
D				0



Element j in a degenerate triple (i, j, k) lies* on the evolutionary path from i to k

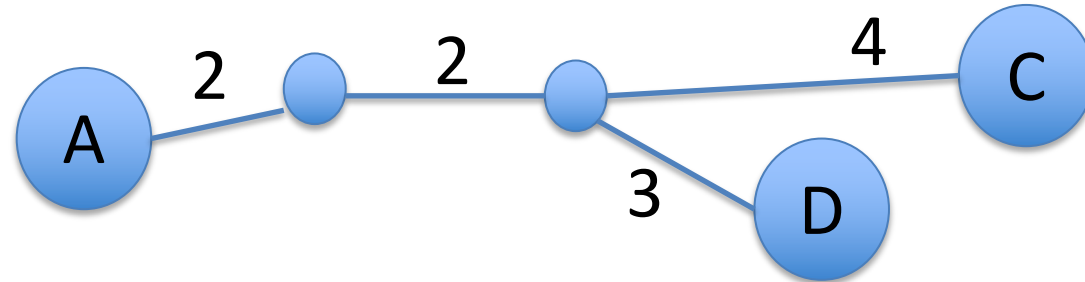
*or is attached to this path by an edge of length 0

Degenerate Triples can be **Removed**

A **degenerate triple** is a set of three distinct elements

$$i, j, k \in [n] \text{ such that } d_{i,j} + d_{j,k} = d_{i,k}$$

	A	C	D
A	0	8	7
C		0	7
D			0

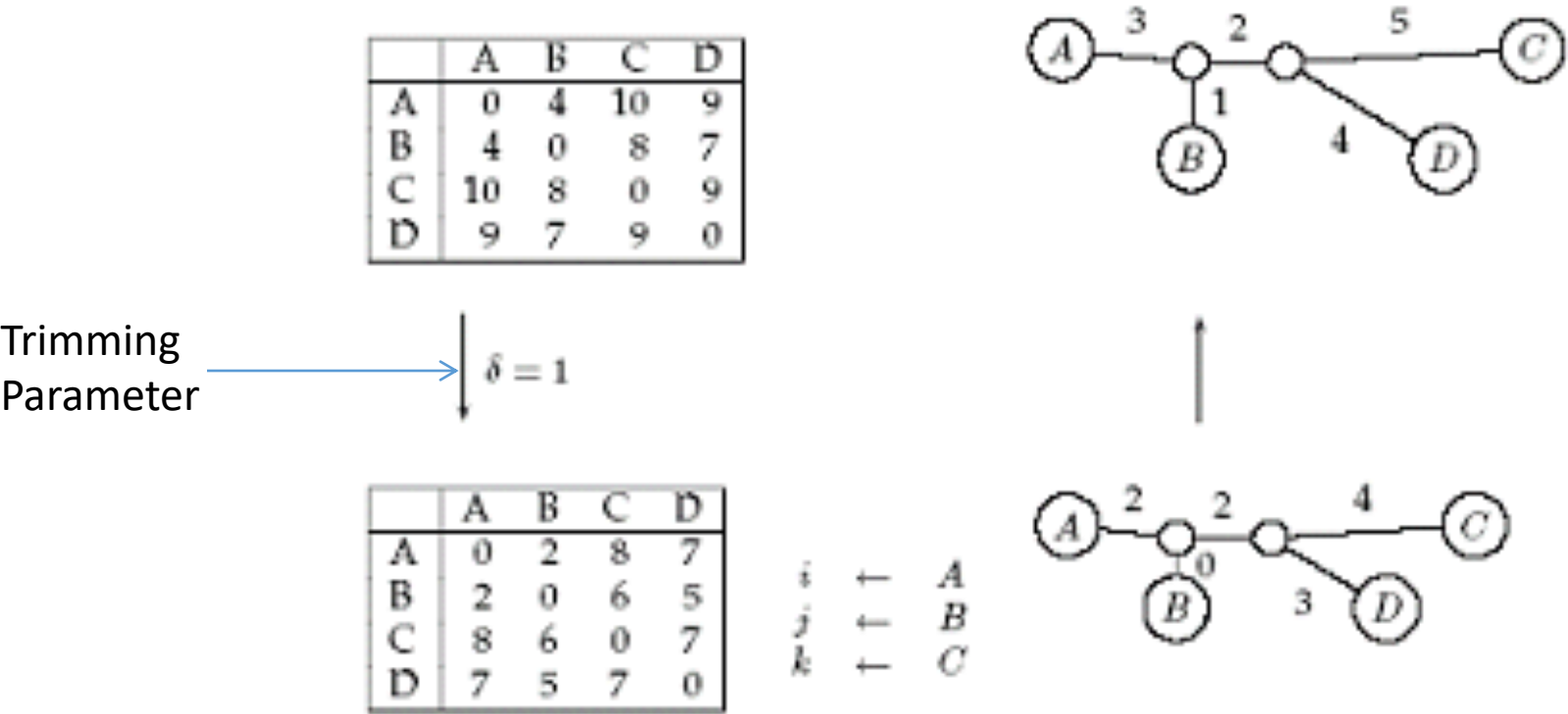


Element j in a degenerate triple (i, j, k) lies* on the evolutionary path from i to k

*or is attached to this path by an edge of length 0

Looking for Degenerate Triples

If distance matrix D does not have a degenerate triple, one can create one by shortening all hanging edges



Decrease entries in matrix D by 2δ

Additive Phylogeny

- If there is no degenerative triple:
 - Reduce all hanging edges by the same amount δ , so that all pairwise distances in the matrix are reduced by 2δ .
- This process will eventually collapse one of the leaves (when δ equals the length of the shortest hanging edge), forming a degenerate triple (i, j, k) and reducing the size of the distance matrix D
- The attachment point for j can be recovered in the reverse transformations by saving $d_{i,j}$ for each collapsed leaf.

	A	B	C	D
A	0	4	10	9
B	4	0	8	7
C	10	8	0	9
D	9	7	9	0

$\delta = 1$

	A	B	C	D
A	0	2	8	7
B	2	0	6	5
C	8	6	0	7
D	7	5	7	0

$i \leftarrow A$
 $j \leftarrow B$
 $k \leftarrow C$

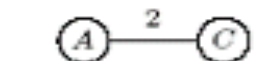
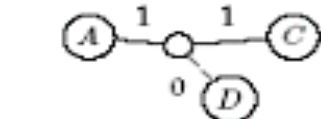
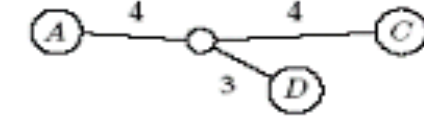
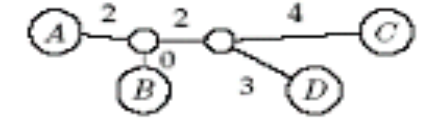
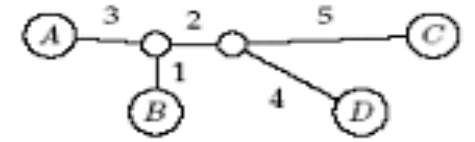
	A	C	D
A	0	8	7
C	8	0	7
D	7	7	0

$\delta = 3$

	A	C	D
A	0	2	1
C	2	0	1
D	1	1	0

$i \leftarrow A$
 $j \leftarrow D$
 $k \leftarrow C$

	A	C
A	0	2
C	2	0



Additive Phylogeny

AdditivePhylogeny(D)

if D is a 2×2 matrix

T = tree of a single edge of length $D_{1,2}$

return T

if D is non-degenerate

Compute trimming parameter δ

Trim(D, δ)

Find a triple i, j, k in D such that $D_{ij} + D_{jk} = D_{ik}$

$x = D_{ij}$

Remove j^{th} row and j^{th} column from D

$T = \text{AdditivePhylogeny}(D)$.

Add a new vertex v to T at distance x from i to k

Add j back to T by creating an edge (v, j) of length 0

for every leaf l in T

if distance from l to v in the tree $\neq D_{lj}$

output "matrix is not additive"

return

Extend all "hanging" edges by length δ

return T

	A	B	C	D
A	0	4	10	9
B	4	0	8	7
C	10	8	0	9
D	9	7	9	0

$\delta = 1$

	A	B	C	D
A	0	2	8	7
B	2	0	6	5
C	8	6	0	7
D	7	5	7	0

$i \leftarrow A$
 $j \leftarrow B$
 $k \leftarrow C$

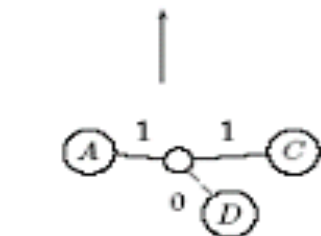
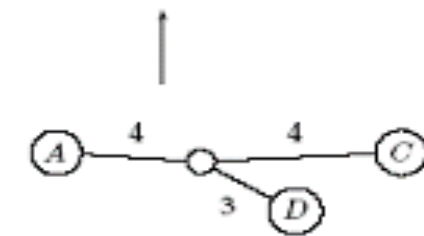
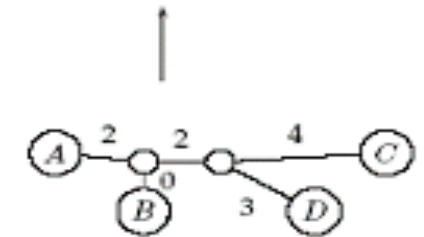
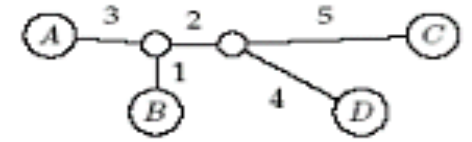
	A	C	D
A	0	8	7
C	8	0	7
D	7	7	0

$\delta = 3$

	A	C	D
A	0	2	1
C	2	0	1
D	1	1	0

$i \leftarrow A$
 $j \leftarrow D$
 $k \leftarrow C$

	A	C
A	0	2
C	2	0




Outline

- Introduction
- Hierarchical clustering
- Additive distance phylogeny
- Four point condition
- Neighbor joining

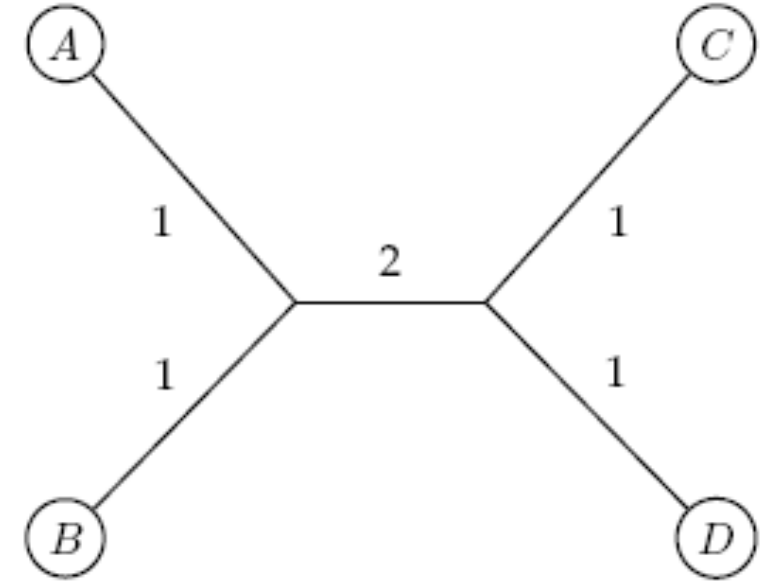
Reading:

- Chapter 10.2 and 10.5-10.8 in Jones and Pevzner

Additive Distance Matrices

Matrix D is  ADDITIVE if there exists a tree T with $d_{ij}(T) = D_{ij}$

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



NON-ADDITIVE
otherwise 

	A	B	C	D
A	0	2	2	2
B	2	0	3	2
C	2	3	0	2
D	2	2	2	0

?

This is a constructive definition

Question: Can we characterize set of additive matrices?

Four Point Condition (Zaretskii 1965, Buneman 1971)

Four point condition of matrix $D = [d_{i,j}]$:

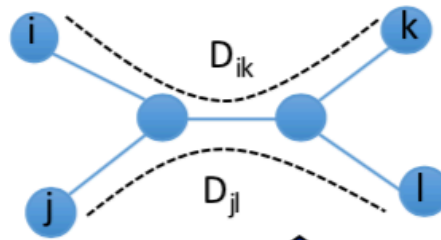
Every four leaves (quartet) can be labeled as (i, j, k, l) such that

$$d_{i,j} + d_{k,l} \leq d_{i,k} + d_{j,l} = d_{i,l} + d_{j,k}$$

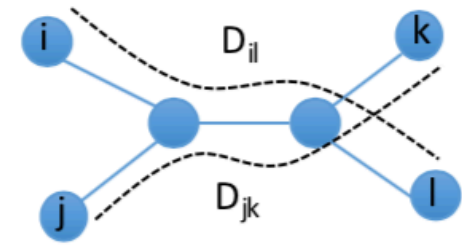
Three sums:

1. $d_{i,j} + d_{k,l}$
2. $d_{i,k} + d_{j,l}$
3. $d_{i,l} + d_{j,k}$

2 and 3 represent the
same number:
(length of all edges)
+ 2 * (length middle
edge)



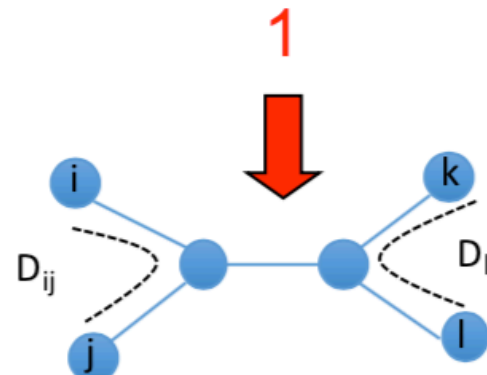
2



3



1 represents a
smaller number:
(length of all edges)
– (length middle
edge)



1



Four Point Condition

Four point condition of matrix $D = [d_{i,j}]$:

Every four leaves (quartet) can be labeled as (i, j, k, l) such that

$$d_{i,j} + d_{k,l} \leq d_{i,k} + d_{j,l} = d_{i,l} + d_{j,k}$$

If two leaves are the same, four point condition is triangle inequality
(e.g. set $l = j$)

Four point condition generalizes triangle inequality and defines a subset of distances, namely additive distances

Four Point Condition: Theorem

Every four leaves (quartet) can be labeled as (i, j, k, l) such that

$$d_{i,j} + d_{k,l} \leq d_{i,k} + d_{j,l} = d_{i,l} + d_{j,k}$$

Theorem: An $n \times n$ matrix D is additive if and only if the four point condition holds for every quartet $(i, j, k, l) \in [n]^4$

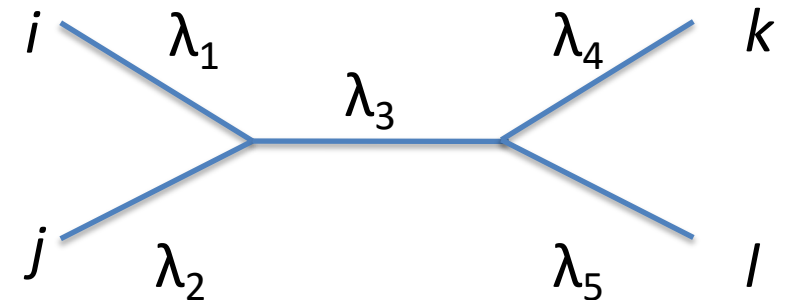
Four Point Condition: Theorem

Every four leaves (quartet) can be labeled as (i, j, k, l) such that

$$d_{i,j} + d_{k,l} \leq d_{i,k} + d_{j,l} = d_{i,l} + d_{j,k}$$

Theorem: An $n \times n$ matrix D is additive if and only if the four point condition holds for every quartet $(i, j, k, l) \in [n]^4$

Proof: (\Rightarrow) Since D is additive, there is a tree T such that $d_{i,j} = d_T(i, j)$ for all $(i, j) \in n^2$. Let (i, j, k, l) be a quartet. Assume w.l.o.g. that i, j and k, l are neighbors. Define λ_m as illustrated.



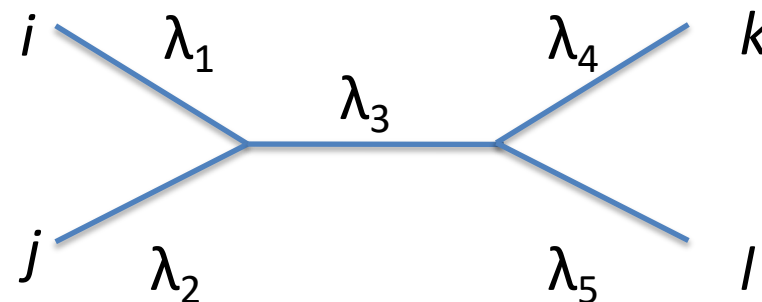
Four Point Condition: Theorem

Every four leaves (quartet) can be labeled as (i, j, k, l) such that

$$d_{i,j} + d_{k,l} \leq d_{i,k} + d_{j,l} = d_{i,l} + d_{j,k}$$

Theorem: An $n \times n$ matrix D is additive if and only if the four point condition holds for every quartet $(i, j, k, l) \in [n]^4$

Proof: (\Rightarrow) Since D is additive, there is a tree T such that $d_{i,j} = d_T(i, j)$ for all $(i, j) \in n^2$. Let (i, j, k, l) be a quartet. Assume w.l.o.g. that i, j and k, l are neighbors. Define λ_m as illustrated.



$$\begin{aligned} d_{i,k} + d_{j,l} &= (\lambda_1 + \lambda_3 + \lambda_4) + (\lambda_2 + \lambda_3 + \lambda_5) = d_{i,l} + d_{j,k} \\ &\geq (\lambda_1 + \lambda_2) + (\lambda_4 + \lambda_5) = d_{i,j} + d_{k,l} \end{aligned}$$

Four Point Condition: Theorem

Every four leaves (quartet) can be labeled as (i, j, k, l) such that

$$d_{i,j} + d_{k,l} \leq d_{i,k} + d_{j,l} = d_{i,l} + d_{j,k}$$

Theorem: An $n \times n$ matrix D is additive if and only if the four point condition holds for every quartet $(i, j, k, l) \in [n]^4$

Proof: (\Leftarrow) Assume four point condition holds. Need an algorithm to construct T . AdditivePhylogeny(T) is one such algorithm*. Neighbor joining is another algorithm.

*we have not proved correctness nor shown how to correct δ

Additive Distance Matrix

Four point condition of matrix $D = [d_{i,j}]$:

Every four leaves (quartet) can be labeled as (i, j, k, l) such that

$$d_{i,j} + d_{k,l} \leq d_{i,k} + d_{j,l} = d_{i,l} + d_{j,k}$$

Theorem: Let D be an $n \times n$ matrix. The following statements are equivalent.

1. Matrix D is additive.
2. There exists a unique tree T (modulo isomorphism) s.t. $d_{i,j} = d_T(i, j)$ for all $(i, j) \in n^2$.
3. Four point condition holds for every quartet $(i, j, k, l) \in [n]^4$.

Outline

- Introduction
- Hierarchical clustering
- Additive distance phylogeny
- Four point condition
- Neighbor joining

Reading:

- Chapter 10.2 and 10.5-10.8 in Jones and Pevzner

Distance Based Phylogeny Problem

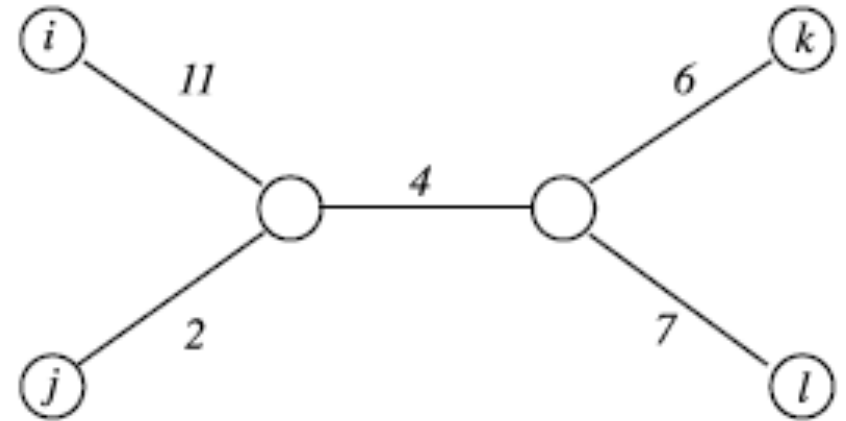
Large Additive Distance Phylogeny Problem:

Given $n \times n$ matrix $D = [d_{i,j}]$, find tree T with n leaves and edge weights such that $\max_{(i,j) \in [n]^2} |d_T(i,j) - d_{i,j}|$ is minimum.

Equivalently, find additive matrix D' closest to input matrix D

Neighbor Joining Algorithm (Saitou and Nei 1987)

- Constructs binary unrooted trees.
- Recall: leaves a and b are neighbors if they have a common parent
- Recall: closest leaves are not necessarily neighbors
- NJ: Find pair of leaves that are “close” to each other but “far” from other leaves

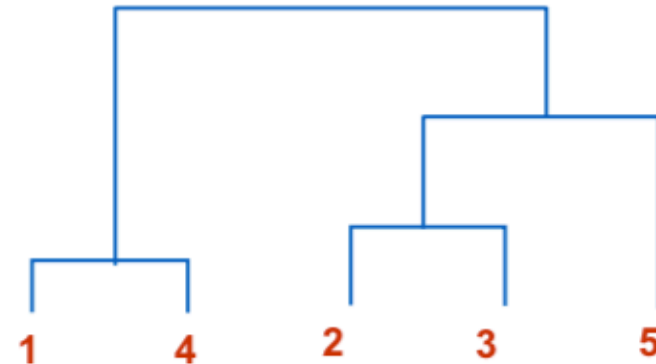
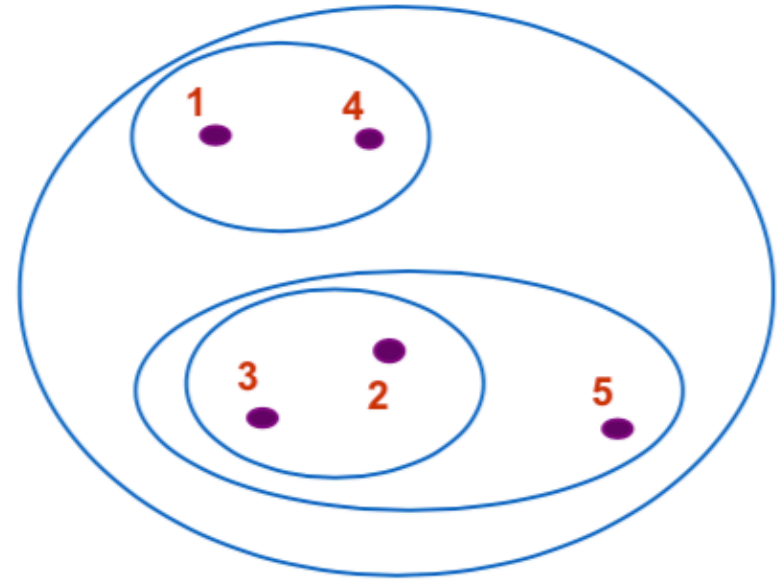


Two advantages: (1) reproduces correct tree for additive matrix, and (2) otherwise gives good approximation of correct tree

Distance Trees as Hierarchical Clustering

Leaves = Data points.

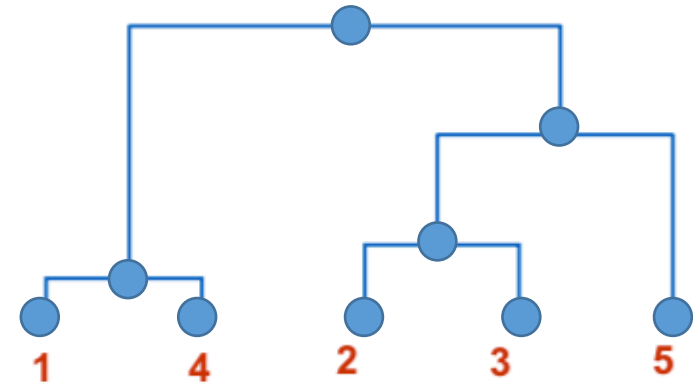
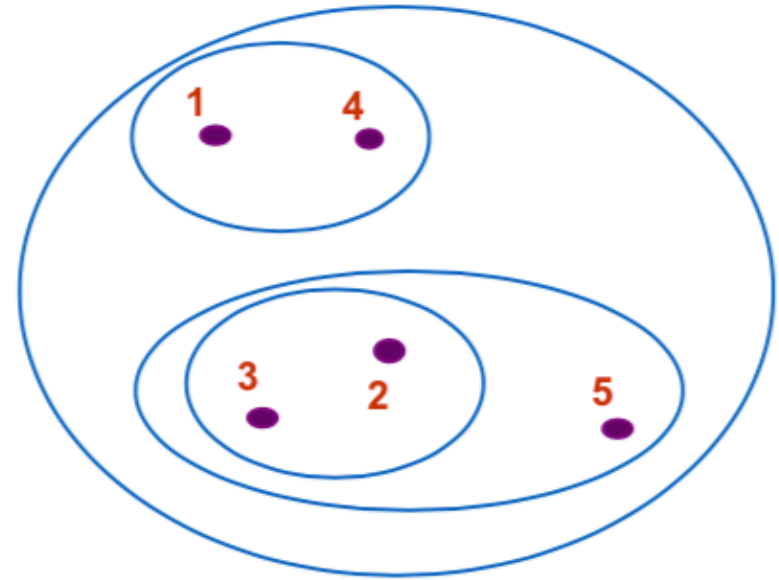
Data points clustered/grouped into hierarchy according to some distance criterion.



Distance Trees as Hierarchical Clustering

Leaves = Data points.

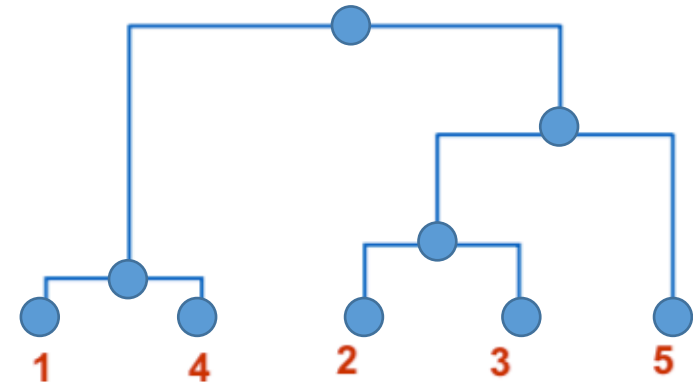
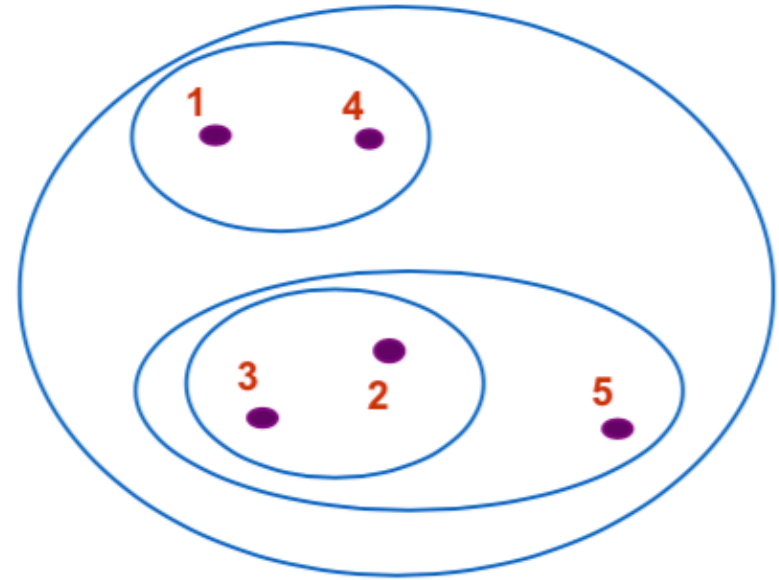
Data points clustered/grouped into hierarchy according to some distance criterion.



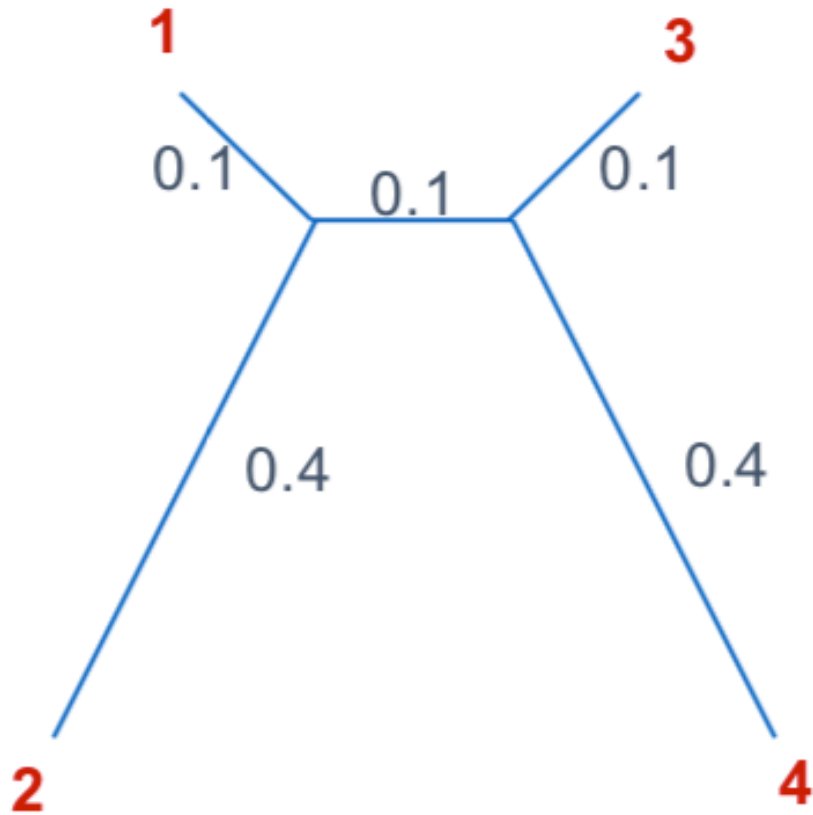
Distance Trees as Hierarchical Clustering

1. Hierarchical Clustering (\mathbf{D} , n)
2. Form n clusters each with one element
3. Construct a graph \mathbf{T} by assigning one vertex to each cluster
4. **while** there is more than one cluster
5. Find the two closest clusters C_1 and C_2
6. Merge C_1 and C_2 into new cluster C with $|C_1| + |C_2|$ elements
7. **Compute distance from C to all other clusters**
8. Add a new vertex \mathbf{C} to \mathbf{T} and connect to vertices C_1 and C_2
9. Remove rows and columns of \mathbf{D} corresponding to C_1 and C_2
10. Add a row and column to \mathbf{D} corresponding to the new cluster \mathbf{C}
11. **return** \mathbf{T}

Selection criterion: distance
between clusters affects
clustering!



Neighbor Joining: Selection Criterion



Let $\mathbf{C} = \{1, \dots, n\}$ be current clusters/leaves.

Define: $u_i = \sum_k D(i, k)$.

Intuitively, u_i measures separation of i from other leaves.

Goal: Minimize $D(i, j)$ and maximize $u_i + u_j$.

Solution: Find pair (i, j) that minimizes:

$$S_D(i, j) = (n - 2) D(i, j) - u_i - u_j$$

Claim: Given additive matrix D .

$S_D(x, y) = \min S_D(i, j)$ if and only if x and y are neighbors in tree T with $d_T = D$.

Neighboring Joining: Algorithm

Initialization:

Form n clusters C_1, C_2, \dots, C_n , one for each leaf node.

Define tree T to be the set of leaf nodes, one per sequence.

Iteration: (D is $m \times m$)

Pick i, j such that $S_D(i, j) = (m - 2) D(i, j) - u_i - u_j$ is minimal.

Merge i and j into new node $[ij]$ in T .

Assign length $\frac{1}{2} (D(i, j) + 1/(m-2) (u_i - u_j))$ to edge $(i, [ij])$

Assign length $\frac{1}{2} (D(i, j) + 1/(m-2) (u_j - u_i))$ to edge $(j, [ij])$

Remove rows and columns from D corresponding to i and j .

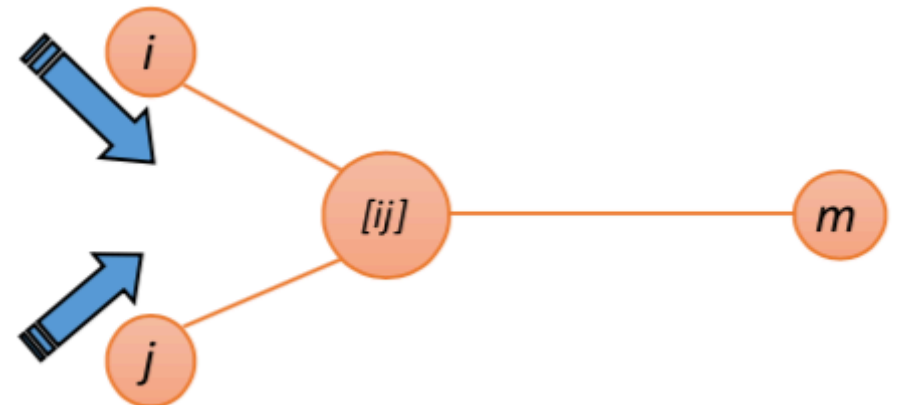
Add row and column to D for new vertex $[ij]$.

Set $D([ij], m) = \frac{1}{2} [D(i, m) + D(j, m) - D(i, j)]$

Termination:

When only one cluster

Question: Does this create rooted or unrooted trees?

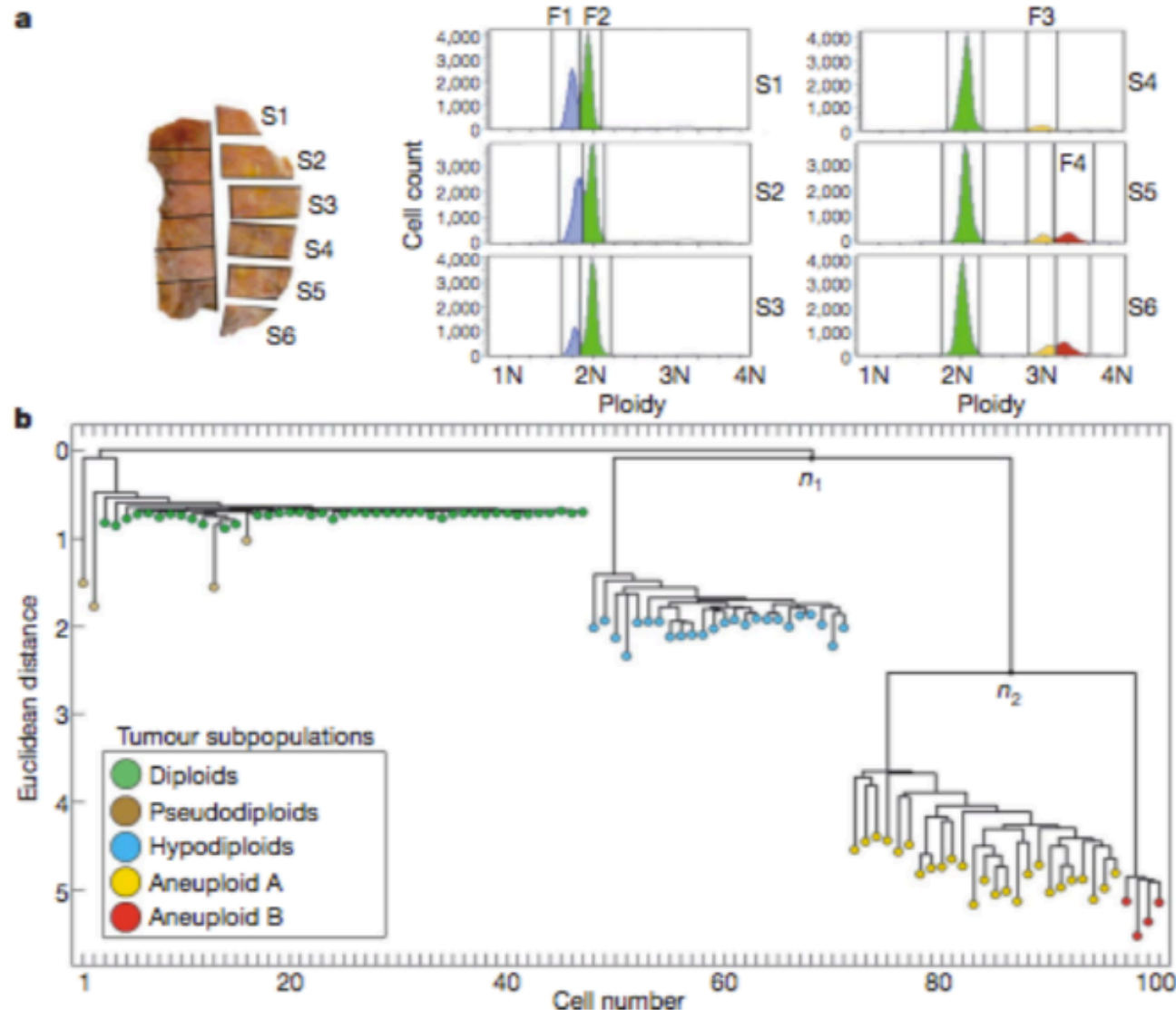


Advantages of Neighbor Joining

Theorem: Let D be an $n \times n$ matrix. If matrix D is additive then neighbor joining produces the unique phylogenetic tree T (modulo isomorphism) such that $d_{i,j} = d_T(i,j)$ for all $(i,j) \in n^2$.

Theorem: Let D be an $n \times n$ matrix. If there exists an additive matrix D' such that $|D - D'|_\infty \leq 0.5$ then neighbor joining applied to D reconstructs the unique tree T (modulo isomorphism) such that $d'_{i,j} = d_T(i,j)$ for all $(i,j) \in n^2$.

Neighbor Joining in Practice



Neighbor Joining tree relating copy number profiles from single cells in a tumor.

[Navin et al, Nature 2011]

Summary

- Introduction
- Hierarchical clustering
- Additive distance phylogeny
- Four point condition
- Neighbor joining

Reading:

- Chapter 10.2 and 10.5-10.8 in Jones and Pevzner