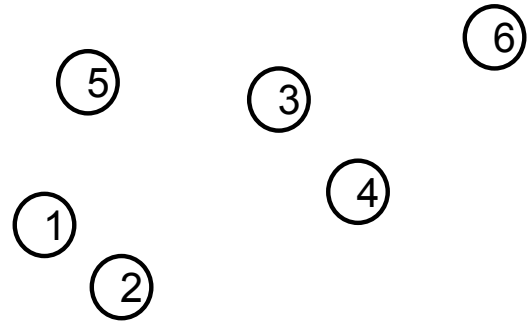
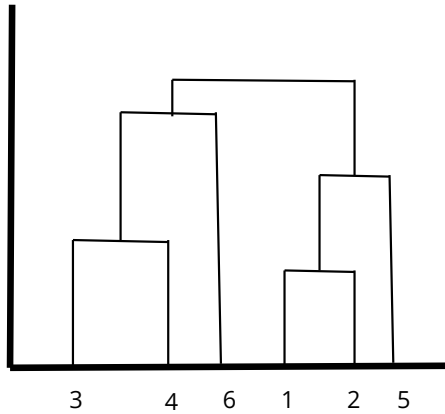

Hierarchical Clustering

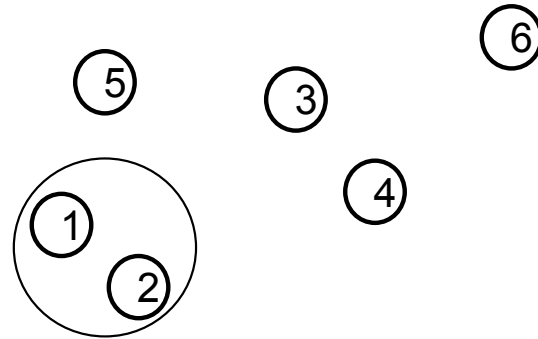
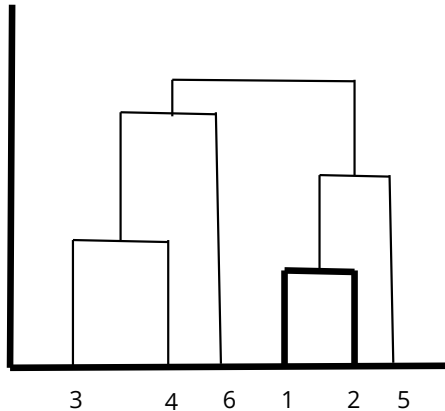
— Boston University CS 506 - Lance Galletti —

Hierarchical Clustering



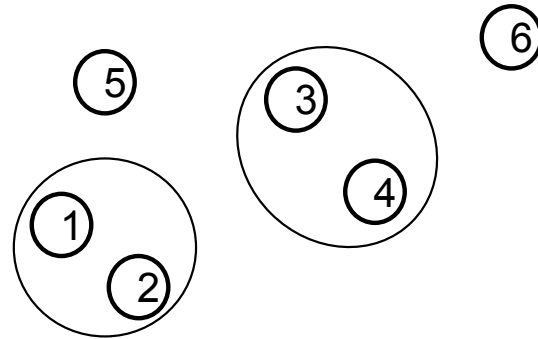
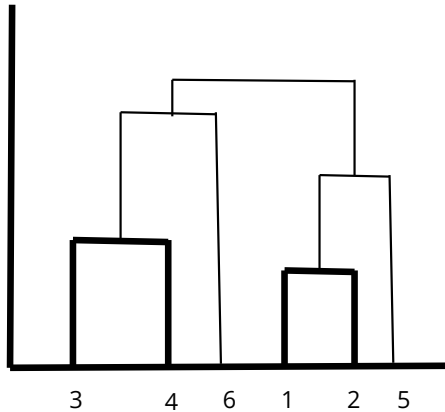
Hierarchical Clustering

At every step, we record which clusters were merged in order to produce a dendrogram:



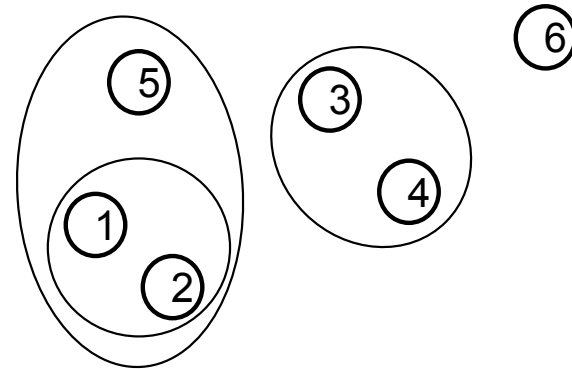
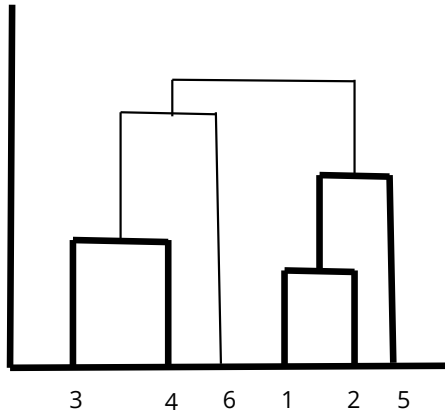
Hierarchical Clustering

At every step, we record which clusters were merged in order to produce a dendrogram:



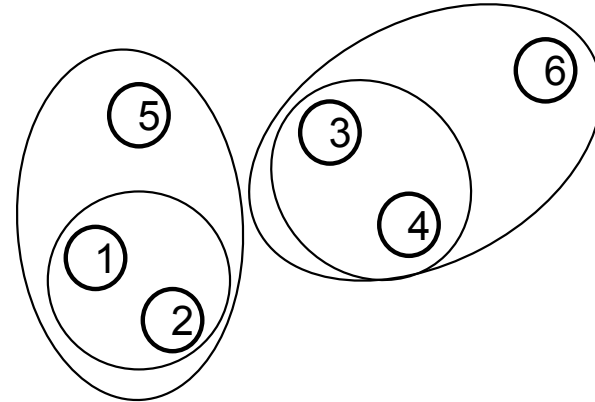
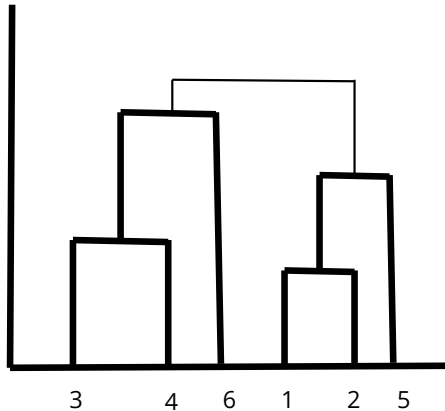
Hierarchical Clustering

At every step, we record which clusters were merged in order to produce a dendrogram:



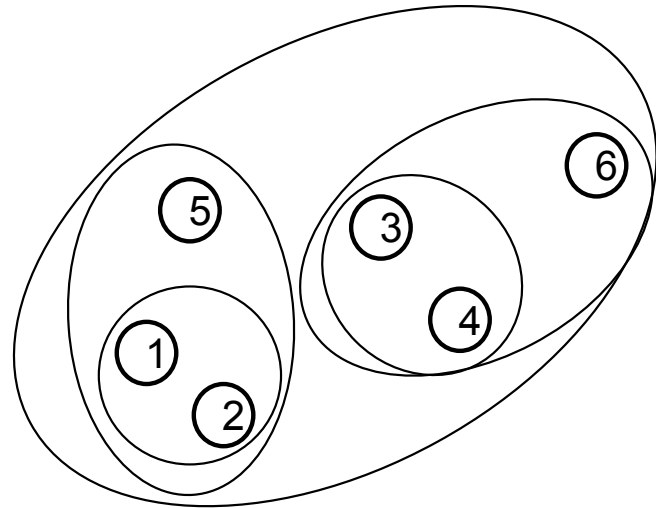
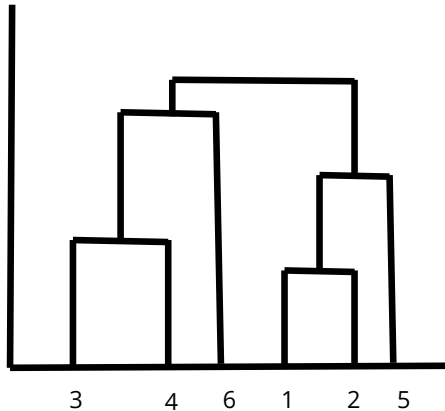
Hierarchical Clustering

At every step, we record which clusters were merged in order to produce a dendrogram:



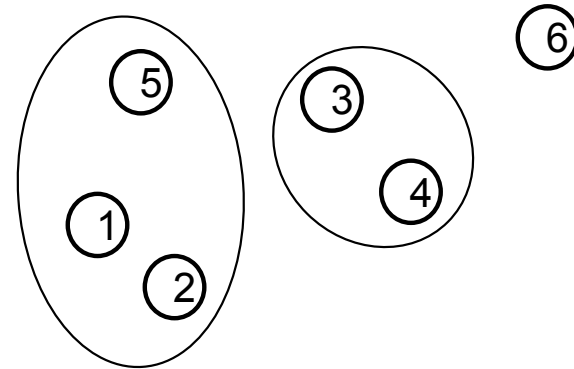
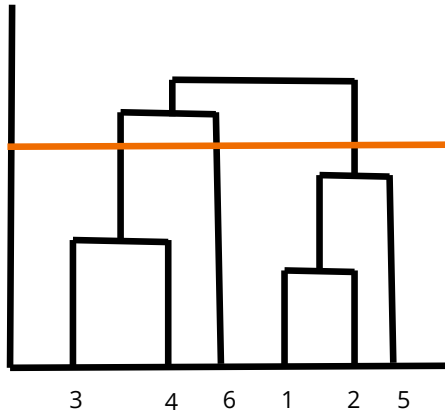
Hierarchical Clustering

At every step, we record which clusters were merged in order to produce a dendrogram:



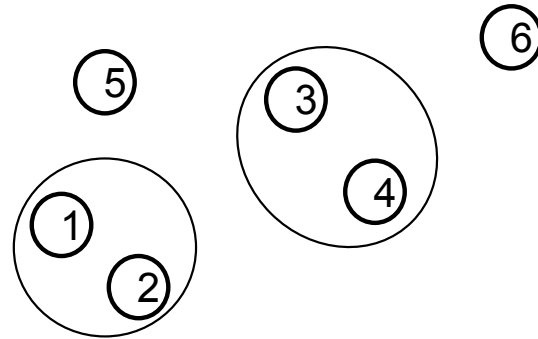
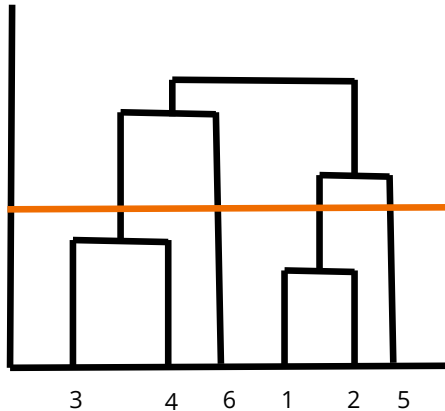
Hierarchical Clustering

We can “cut” the dendrogram at any threshold to produce any number of clusters



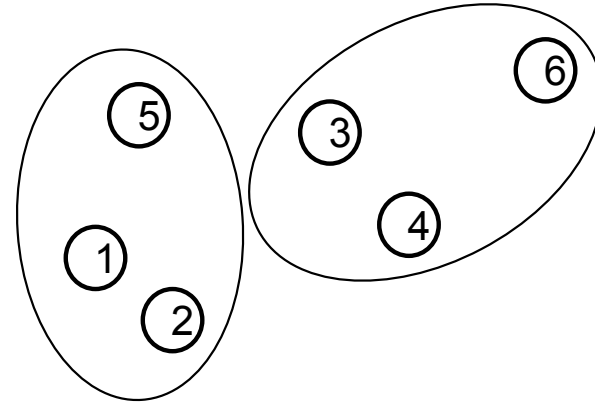
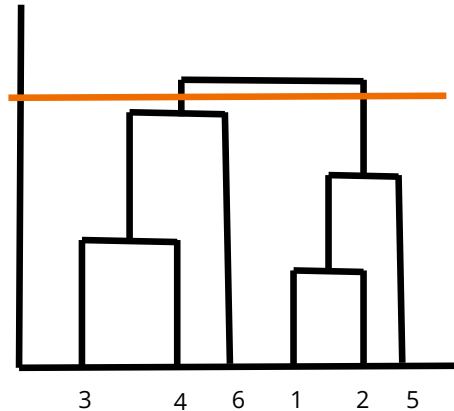
Hierarchical Clustering

We can “cut” the dendrogram at any threshold to produce any number of clusters



Hierarchical Clustering

We can “cut” the dendrogram at any threshold to produce any number of clusters



Example

<https://www.catalogueoflife.org/>

Hierarchical Clustering

Two types of hierarchical clustering:

Agglomerative: *“from the ground up”*

1. Start with every point in its own cluster
2. At each step, merge the two closest clusters
3. Stop when every point is in the same cluster

Divisive: *“every point in one big cluster → divide the cluster”*

1. Start with every point in the same cluster
2. At each step, split until every point is in its own cluster

Agglomerative Clustering Algorithm

1. Let each point in the dataset be in its own cluster
2. Compute the distance between all pairs of clusters
3. Merge the two closest clusters
4. Repeat 3 & 4 until all points are in the same cluster

at every step we merge the two closest clusters
repeat until every point is in one cluster

Hierarchical Clustering

Can we implement this? Are we missing anything?

↳ how do we compute distance between clusters

How would you define distance between clusters?

- ↳ distance between two centroids of the clusters
- ↳ pairwise distances between points and average them out

Hierarchical Clustering - Distance Functions

Let's first define:

Distance between points: $d(p_1, p_2)$

Distance between clusters: $D(C_1, C_2)$

Single-Link Distance

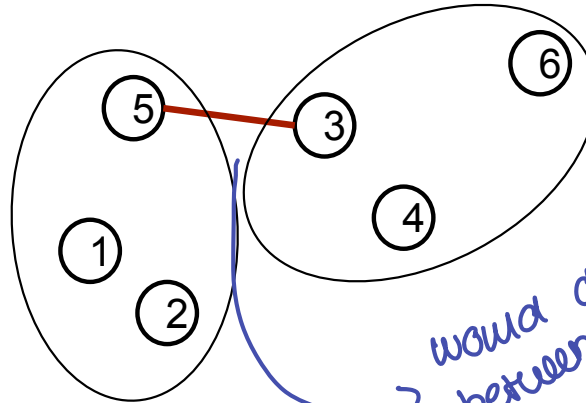
Is the **minimum** of all pairwise distances between a point from one cluster and a point from the other cluster.

↓ distance between p_1 and p_2

$$D_{SL}(C_1, C_2) = \min \{d(p_1, p_2) \mid p_1 \in C_1, p_2 \in C_2\}$$

“shortest distance between two points”

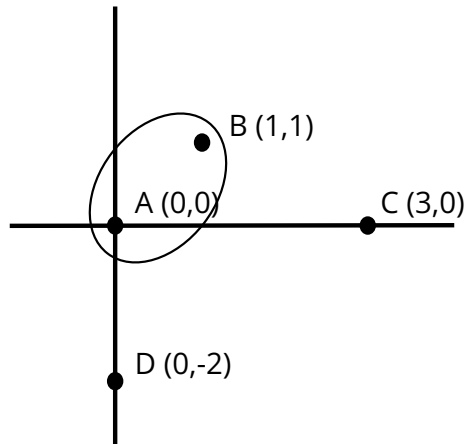
distance between clusters



would define the distance between two clusters

Depends on choice of d

Q: Single-Link Distance



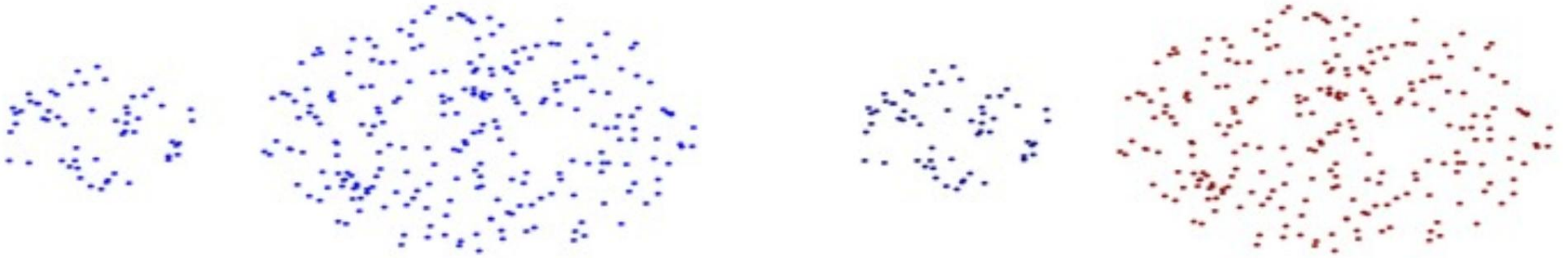
$$A \rightarrow C = 2$$
$$B \rightarrow C = \sqrt{5}$$

Is C or D closer to {A, B} ?

$$D(c(C), c(A, B)) = \sqrt{5}$$

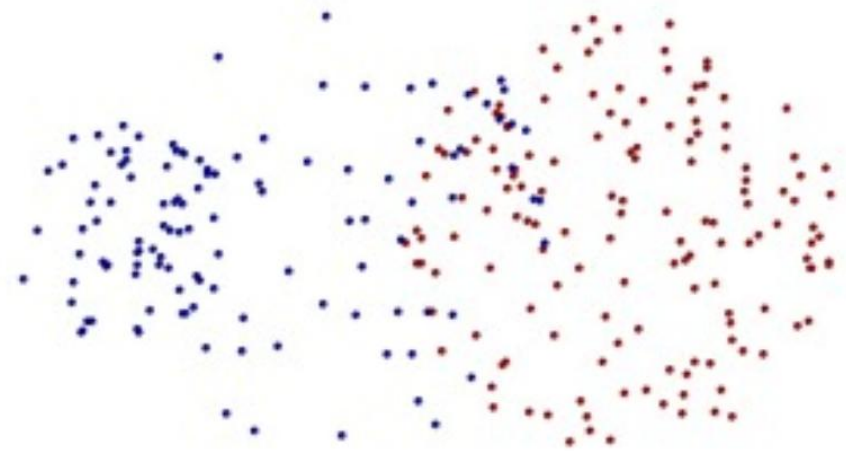
$$A \rightarrow D = 2$$
$$B \rightarrow D = \sqrt{10}$$

Single-Link Distance



Can handle clusters of different sizes

Single-Link Distance



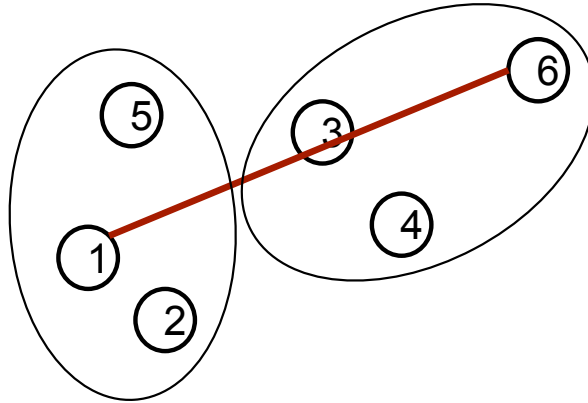
But... Sensitive to noise points
Tends to create elongated clusters

Complete-Link Distance

Is the **maximum** of all pairwise distances between a point from one cluster and a point from the other cluster.

$$D_{CL}(C_1, C_2) = \max \{d(p_1, p_2) \mid p_1 \in C_1, p_2 \in C_2\}$$

max distance between points



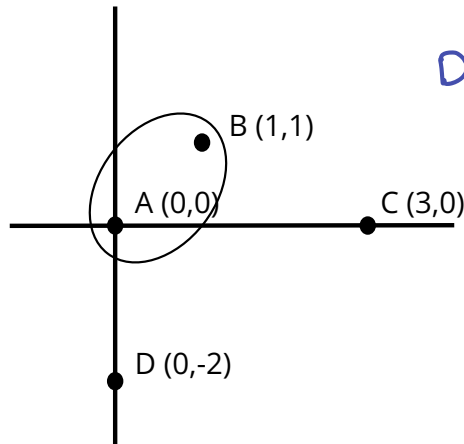
Q: Complete-Link Distance

using maximum

$$C \rightarrow A, B: 3$$

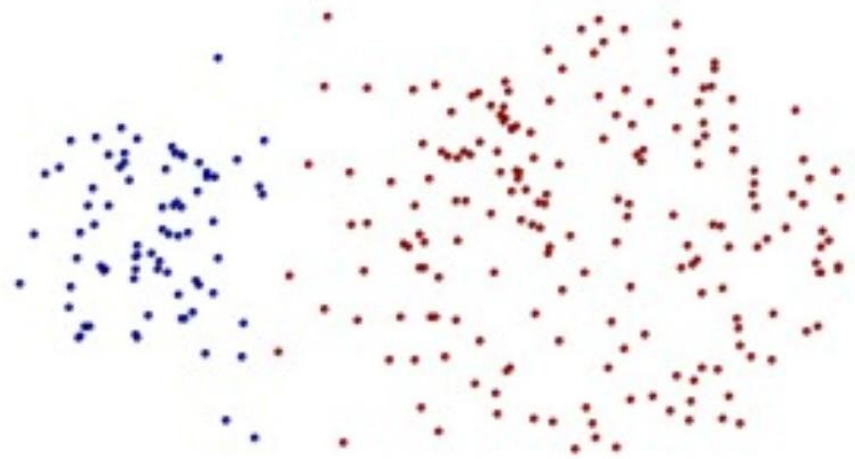
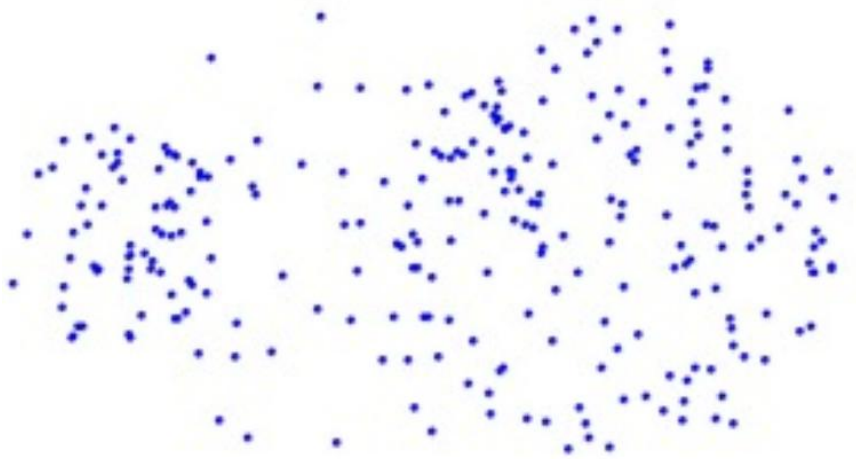
Is C or D closer to $\{A, B\}$?

$$D \rightarrow A, B: \sqrt{10}$$



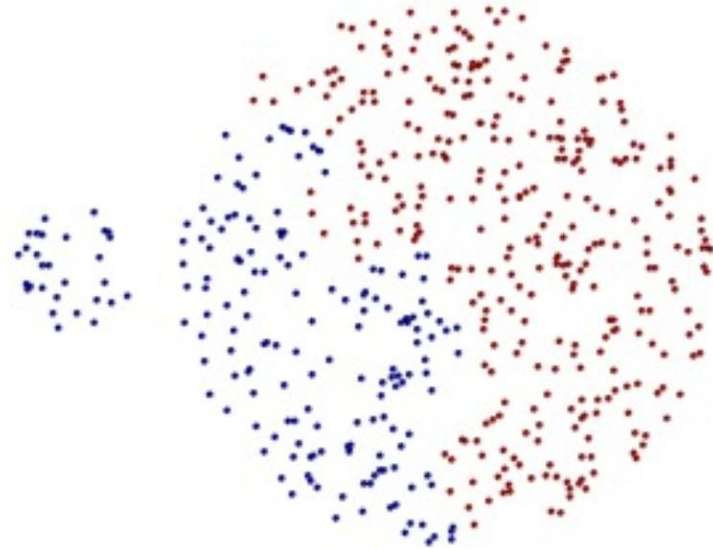
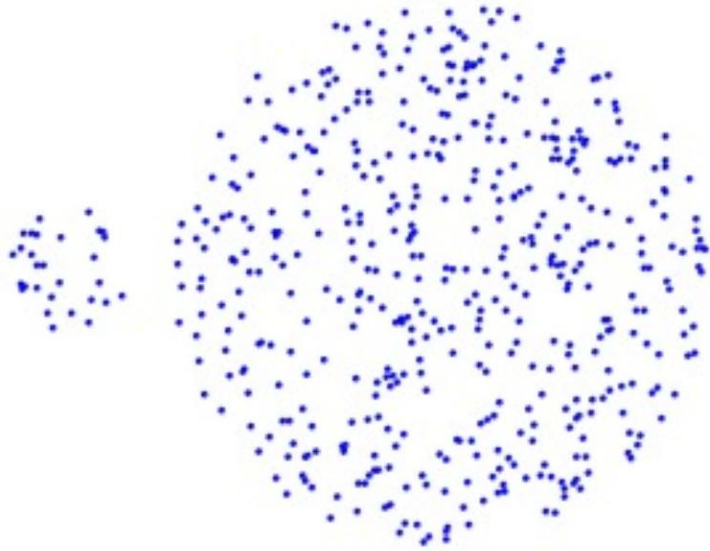
C is closer

Complete-Link Distance



Less susceptible to noise
Creates more balanced (equal diameter) clusters

Complete-Link Distance



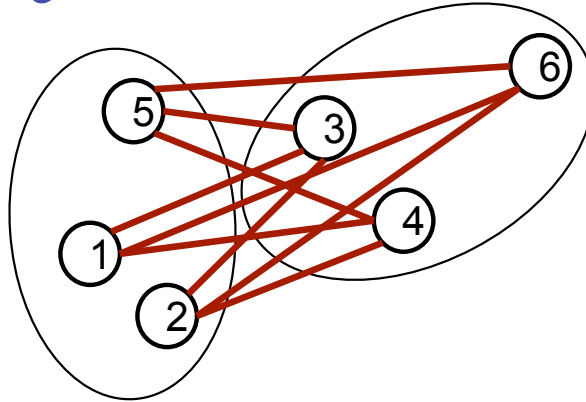
But... Tends to split up large clusters.
All clusters tend to have the same diameter

Average-Link Distance

Is the **average** of all pairwise distances between a point from one cluster and a point from the other cluster.

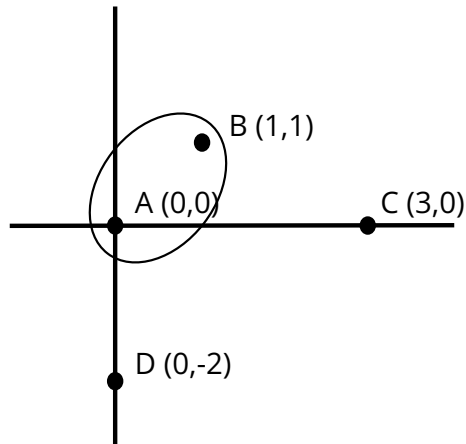
$$D_{AL}(C_1, C_2) = \frac{1}{|C_1| \cdot |C_2|} \sum_{p_1 \in C_1, p_2 \in C_2} d(p_1, p_2)$$

cardinality



Q: Average-Link Distance

Is C or D closer to $\{A, B\}$?



Average-Link Distance

Less susceptible to noise and outliers.

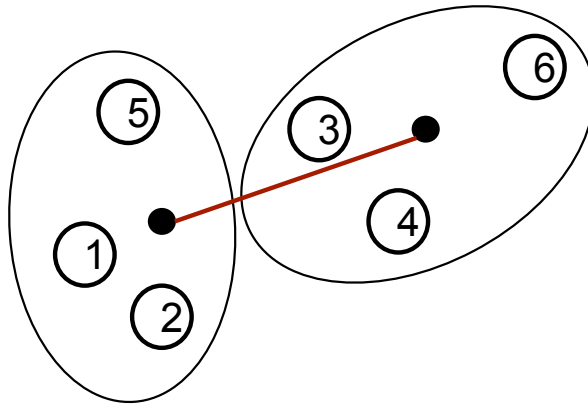
But... Tends to be biased toward globular clusters

Centroid Distance

The distance between the centroids of clusters.

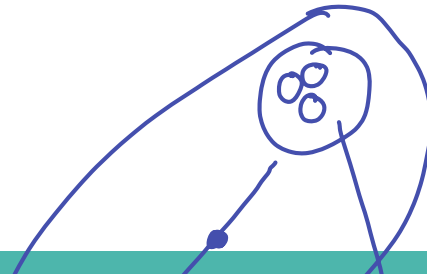
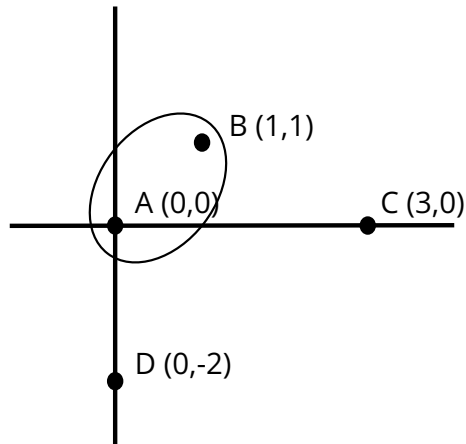
$$D_C(C_1, C_2) = d(\mu_1, \mu_2)$$

↑ centroid₁ ↑ centroid₂



Q: Centroid Distance

Is C or D closer to $\{A, B\}$?



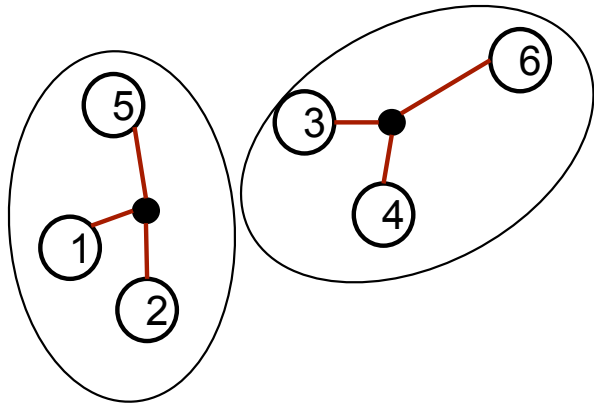
Ward's Distance

↳ used most in practice

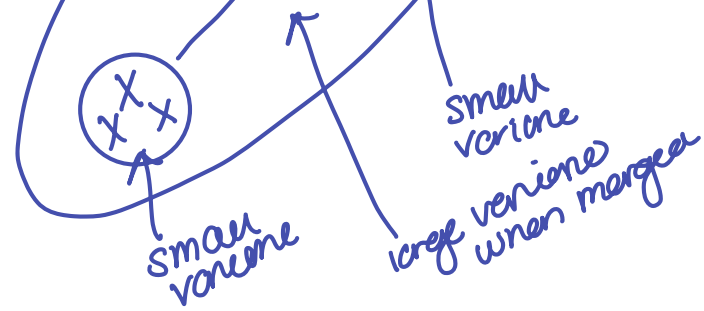
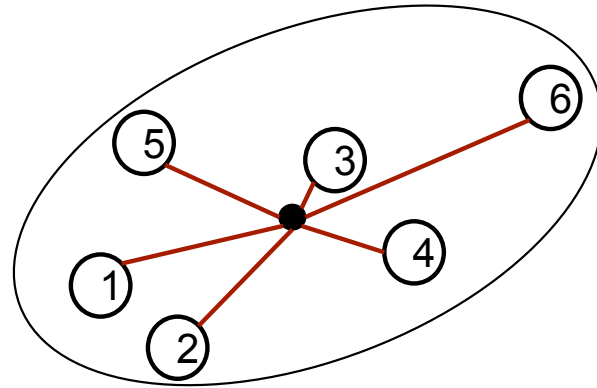
Is the difference between the spread / variance of points in the merged cluster and the unmerged clusters.

$$D_{WD}(C_1, C_2) = \sum_{p \in C_{12}} d(p, \mu_{12}) - \sum_{p_1 \in C_1} d(p_1, \mu_1) - \sum_{p_2 \in C_2} d(p_2, \mu_2)$$

gives notion of variance (pointing to the first sum)
in all clusters (pointing to the last sum)
merge / union of C_1, C_2 (pointing to the first sum)

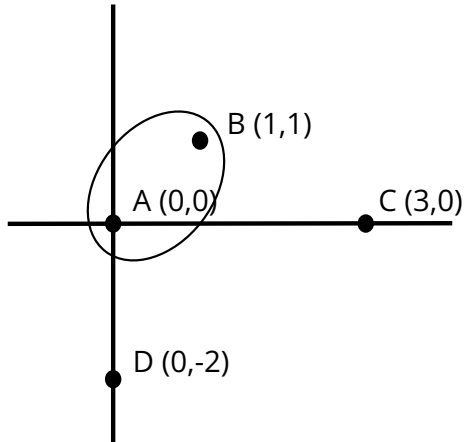


VS



Q: Ward's Distance

Is C or D closer to $\{A, B\}$?



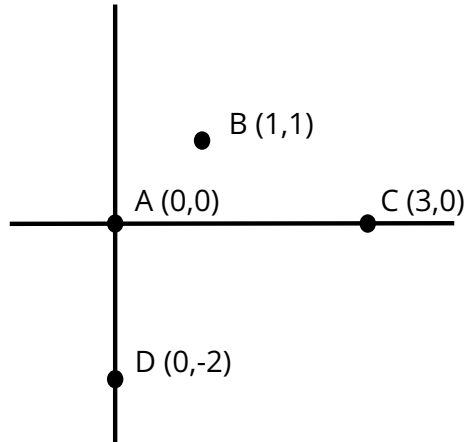
Agglomerative Clustering Algorithm

1. Let each point in the dataset be in its own cluster
2. Compute the distance between all pairs of clusters
3. Merge the two closest clusters
4. Repeat 3 & 4 until all points are in the same cluster

Example

d = Euclidean

D = Single-Link



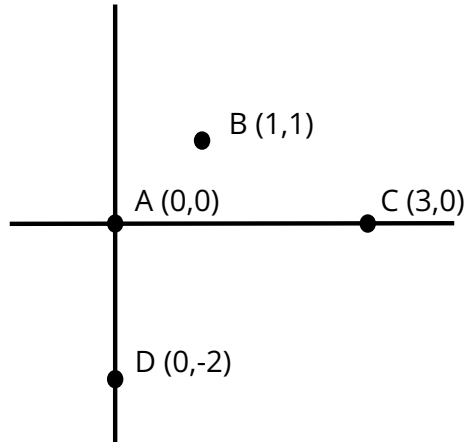
Distance Matrix

	A	B	C	D
A				
B				
C				
D				

Example

d = Euclidean

D = Single-Link



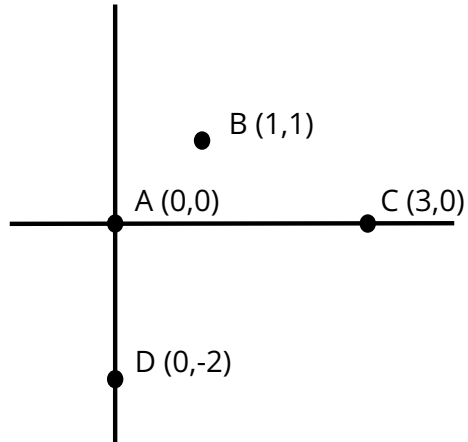
Distance Matrix

	A	B	C	D
A	0			
B		0		
C			0	
D				0

Example

d = Euclidean

D = Single-Link



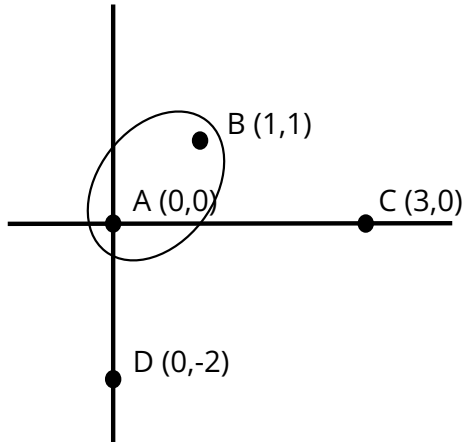
Distance Matrix

	A	B	C	D
A	0	$\sqrt{2}$	3	2
B	$\sqrt{2}$	0	$\sqrt{5}$	$\sqrt{10}$
C	3	$\sqrt{5}$	0	$\sqrt{13}$
D	2	$\sqrt{10}$	$\sqrt{13}$	0

Example

d = Euclidean

D = Single-Link



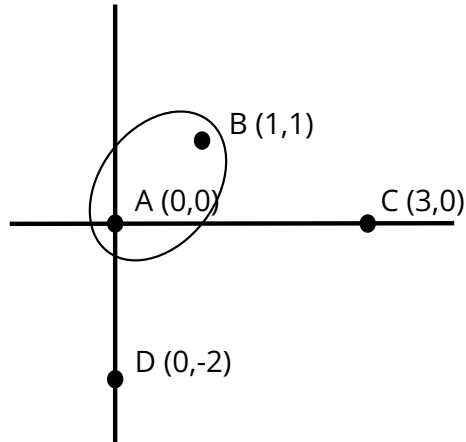
Distance Matrix

	A	B	C	D
A	0	$\sqrt{2}$	3	2
B	$\sqrt{2}$	0	$\sqrt{5}$	$\sqrt{10}$
C	3	$\sqrt{5}$	0	$\sqrt{13}$
D	2	$\sqrt{10}$	$\sqrt{13}$	0

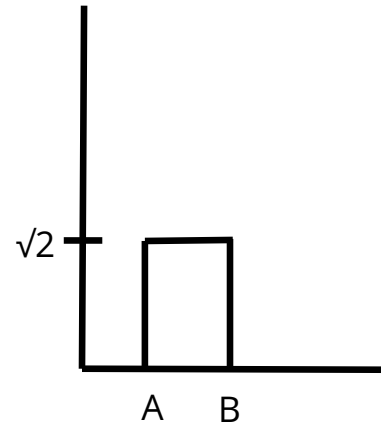
Example

d = Euclidean

D = Single-Link



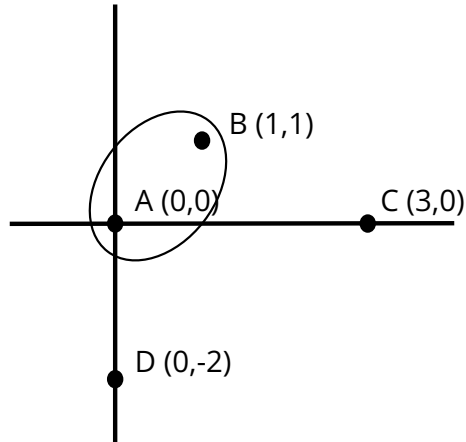
Dendrogram



Example

d = Euclidean

D = Single-Link



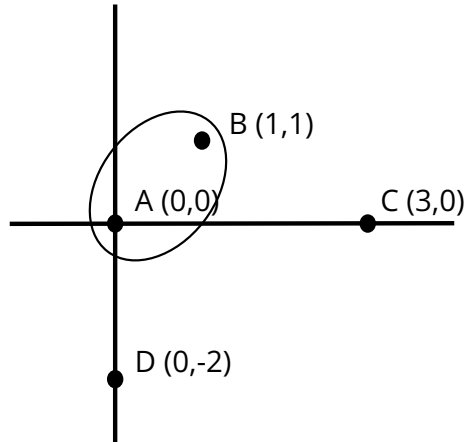
Distance Matrix

	A & B	C	D
A & B	0		
C		0	$\sqrt{13}$
D		$\sqrt{13}$	0

Example

d = Euclidean

D = Single-Link



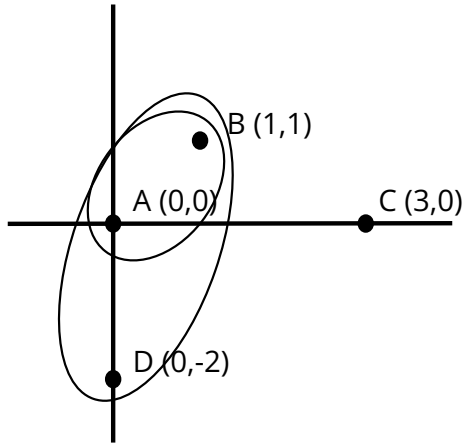
Distance Matrix

	A & B	C	D
A & B	0	$\sqrt{5}$	2
C	$\sqrt{5}$	0	$\sqrt{13}$
D	2	$\sqrt{13}$	0

Example

d = Euclidean

D = Single-Link



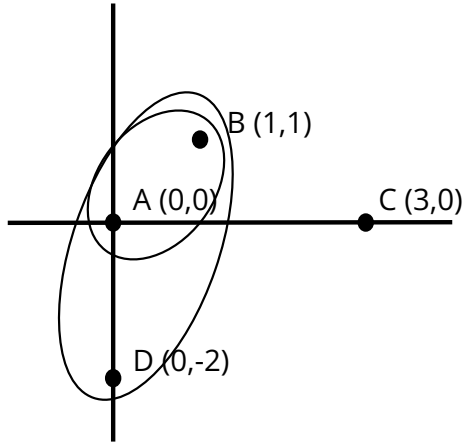
Distance Matrix

	A & B	C	D
A & B	0	$\sqrt{5}$	2
C	$\sqrt{5}$	0	$\sqrt{13}$
D	2	$\sqrt{13}$	0

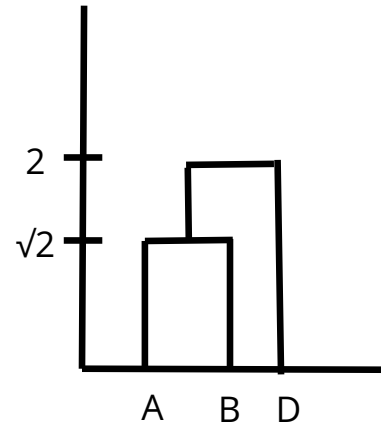
Example

d = Euclidean

D = Single-Link



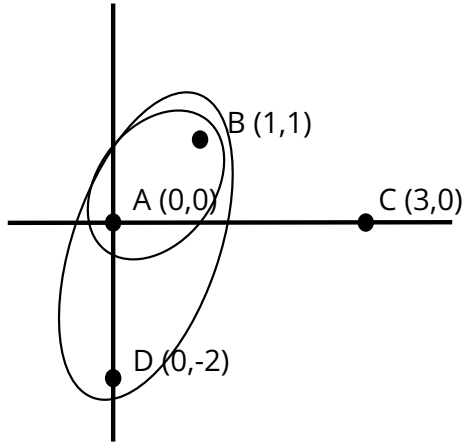
Dendrogram



Example

d = Euclidean

D = Single-Link



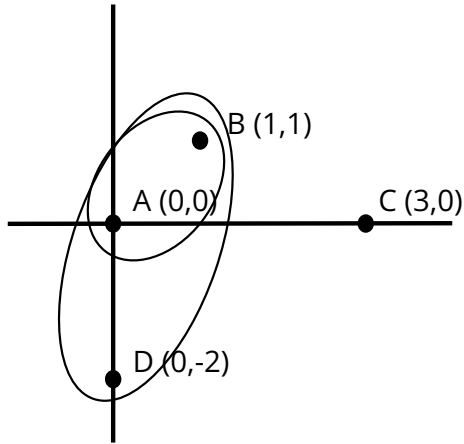
Distance Matrix

	A & B & D	C
A & B & D	0	
C		0

Example

d = Euclidean

D = Single-Link



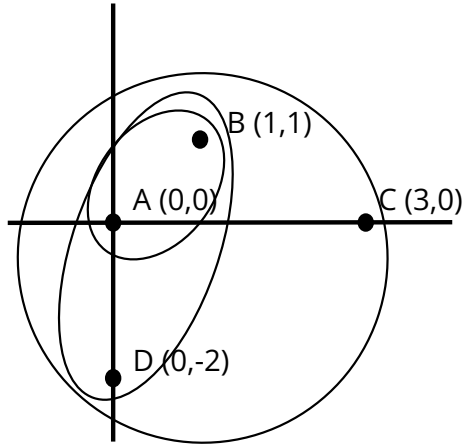
Distance Matrix

	A & B & D	C
A & B & D	0	$\sqrt{5}$
C	$\sqrt{5}$	0

Example

d = Euclidean

D = Single-Link



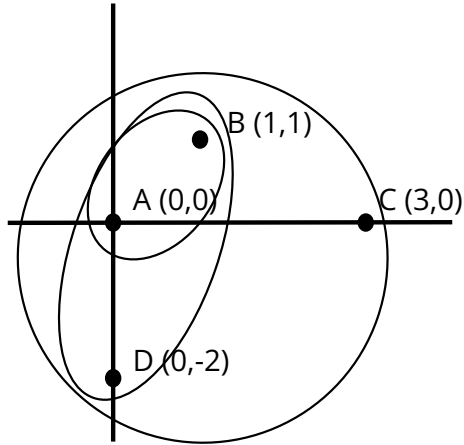
Distance Matrix

	A & B & D	C
A & B & D	0	$\sqrt{5}$
C	$\sqrt{5}$	0

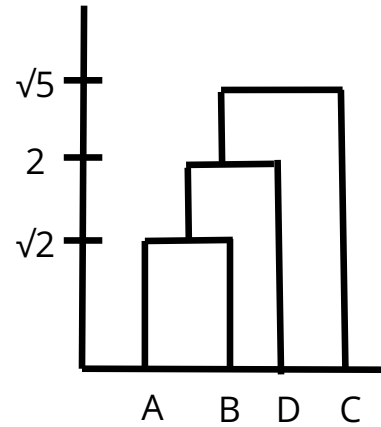
Example

d = Euclidean

D = Single-Link



Dendrogram



Hierarchical Clustering

Finding the threshold with which to cut the dendrogram requires exploration and tuning. But in general hierarchical clustering is used to expose a hierarchy in the data (ex: finding/defining species via DNA similarity).

To capture the difference between clusterings you can use a cost function, or methods that we will discuss later when we look at clustering aggregation.