



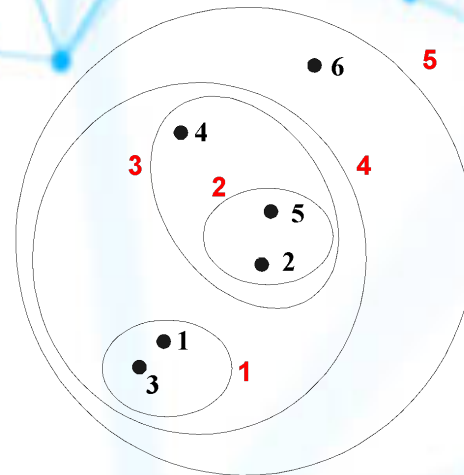
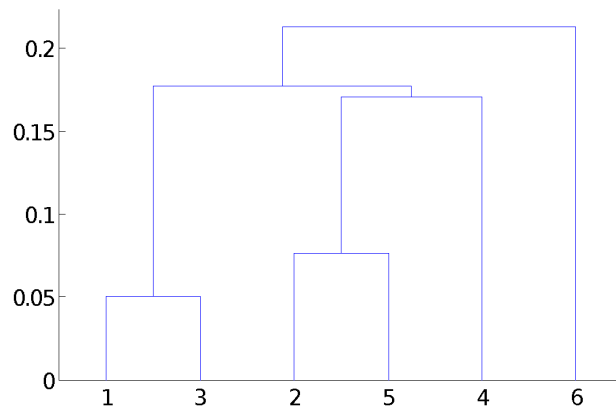
Fundamentals of Big Data Analytics

Lecture 7-8 Hierarchical Clustering

Dr. Iqra Safder
Assistant Professor
FAST NUCES, Lahore

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



Strengths of Hierarchical Clustering

- Do not have to assume any particular number of clusters
 - Any desired number of clusters can be obtained by ‘cutting’ the dendrogram at the proper level
- They may correspond to meaningful taxonomies
 - Example in biological sciences (e.g., animal kingdom, phylogeny reconstruction, ...)

Hierarchical Clustering

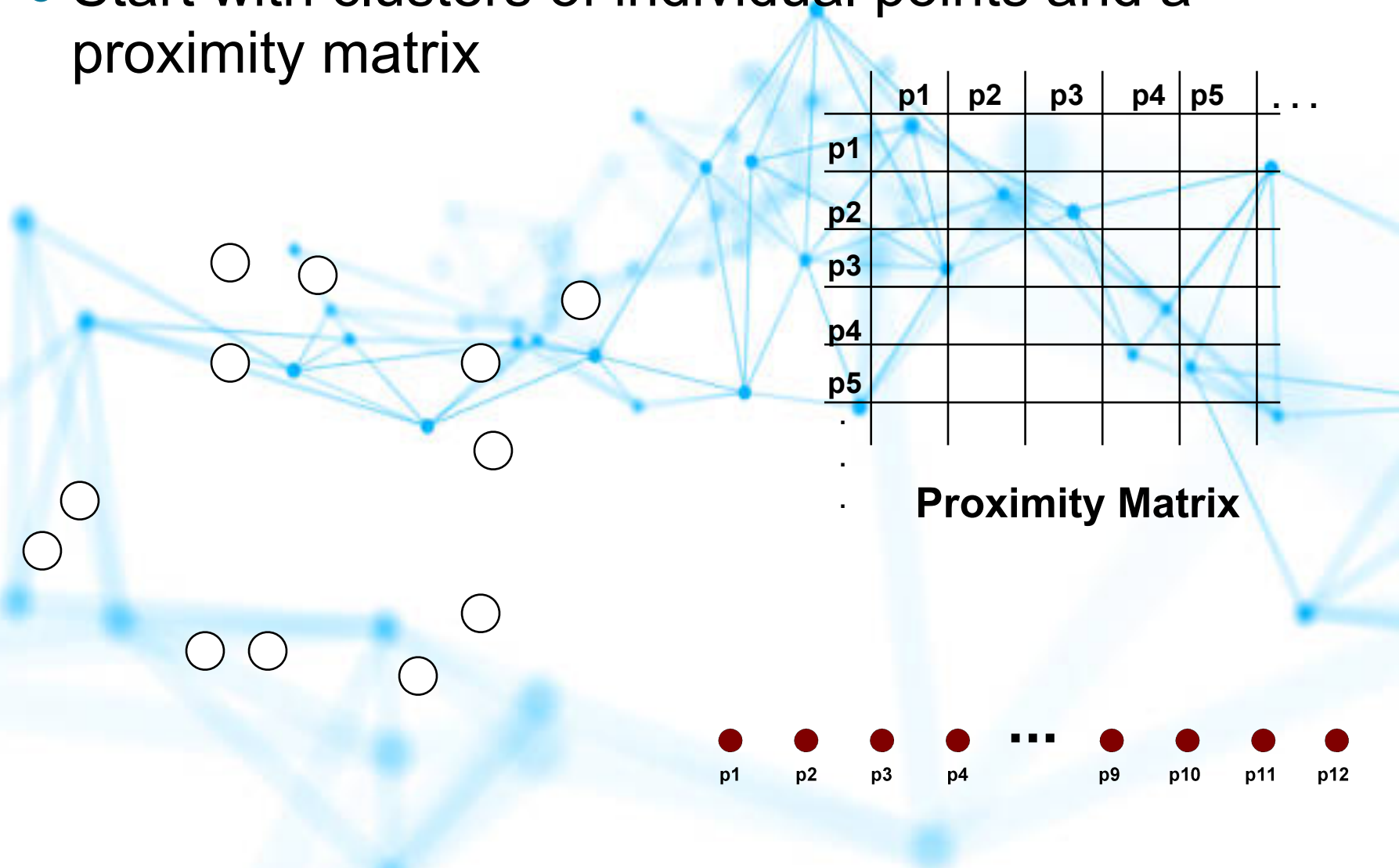
- Two main types of hierarchical clustering
 - Agglomerative:
 - ◆ Start with the points as individual clusters
 - ◆ At each step, merge the closest pair of clusters until only one cluster (or k clusters) left
 - Divisive:
 - ◆ Start with one, all-inclusive cluster
 - ◆ At each step, split a cluster until each cluster contains a point (or there are k clusters)
- Traditional hierarchical algorithms use a similarity or distance matrix
 - Merge or split one cluster at a time

Agglomerative Clustering Algorithm

- More popular hierarchical clustering technique
- Basic algorithm is straightforward
 1. Compute the proximity matrix
 2. Let each data point be a cluster
 3. **Repeat**
 4. Merge the two closest clusters
 5. Update the proximity matrix
 6. **Until** only a single cluster remains
- Key operation is the computation of the proximity of two clusters
 - Different approaches to defining the distance between clusters distinguish the different algorithms

Starting Situation

- Start with clusters of individual points and a proximity matrix



The background of the slide features a network graph with blue nodes and edges, representing a complex structure. To the left of the graph, there are several white circles, some of which are connected to the blue nodes by thin lines. The graph itself is composed of many blue nodes connected by a dense web of edges, with some nodes having more connections than others.

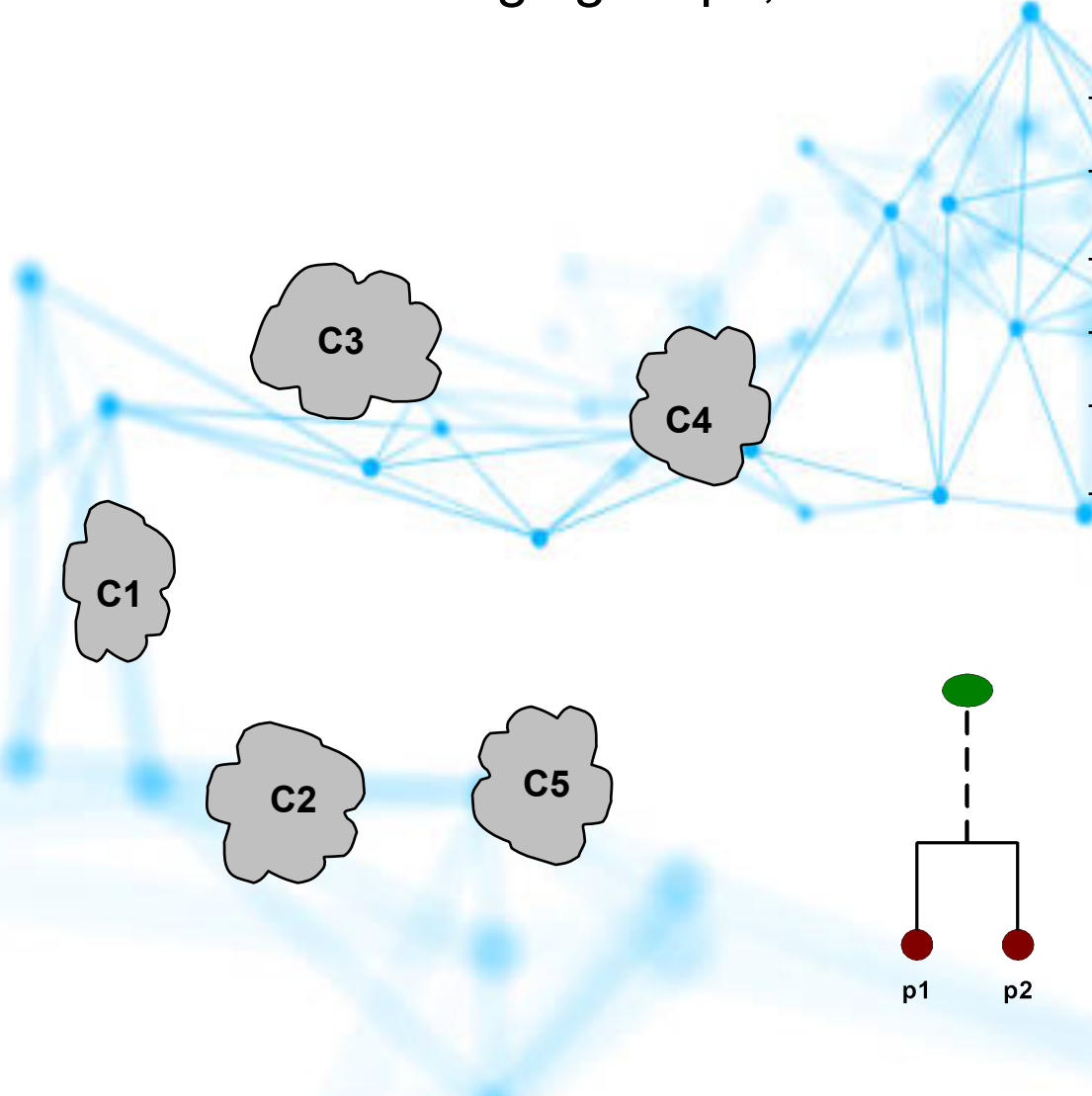
	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						
.						

Proximity Matrix

p1 p2 p3 p4 ... p9 p10 p11 p12

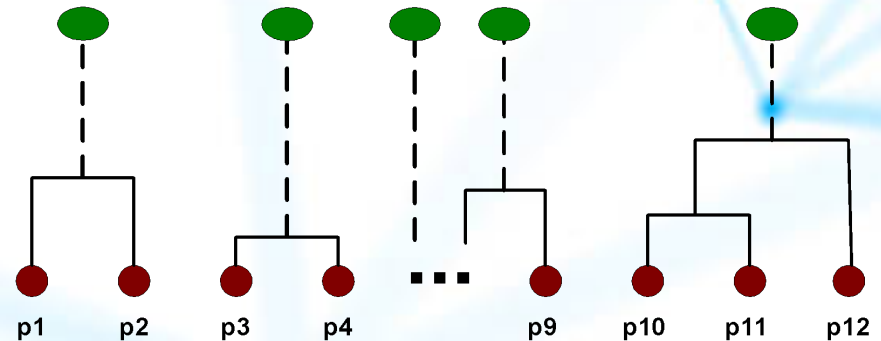
Intermediate Situation

- After some merging steps, we have some clusters



	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Proximity Matrix

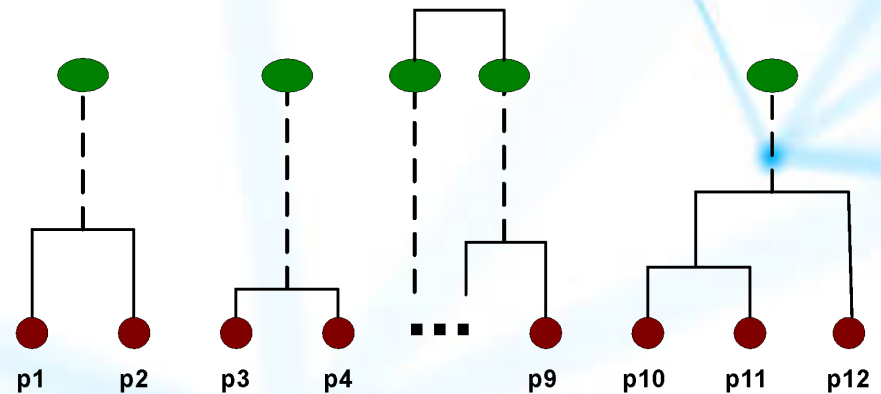


Intermediate Situation

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

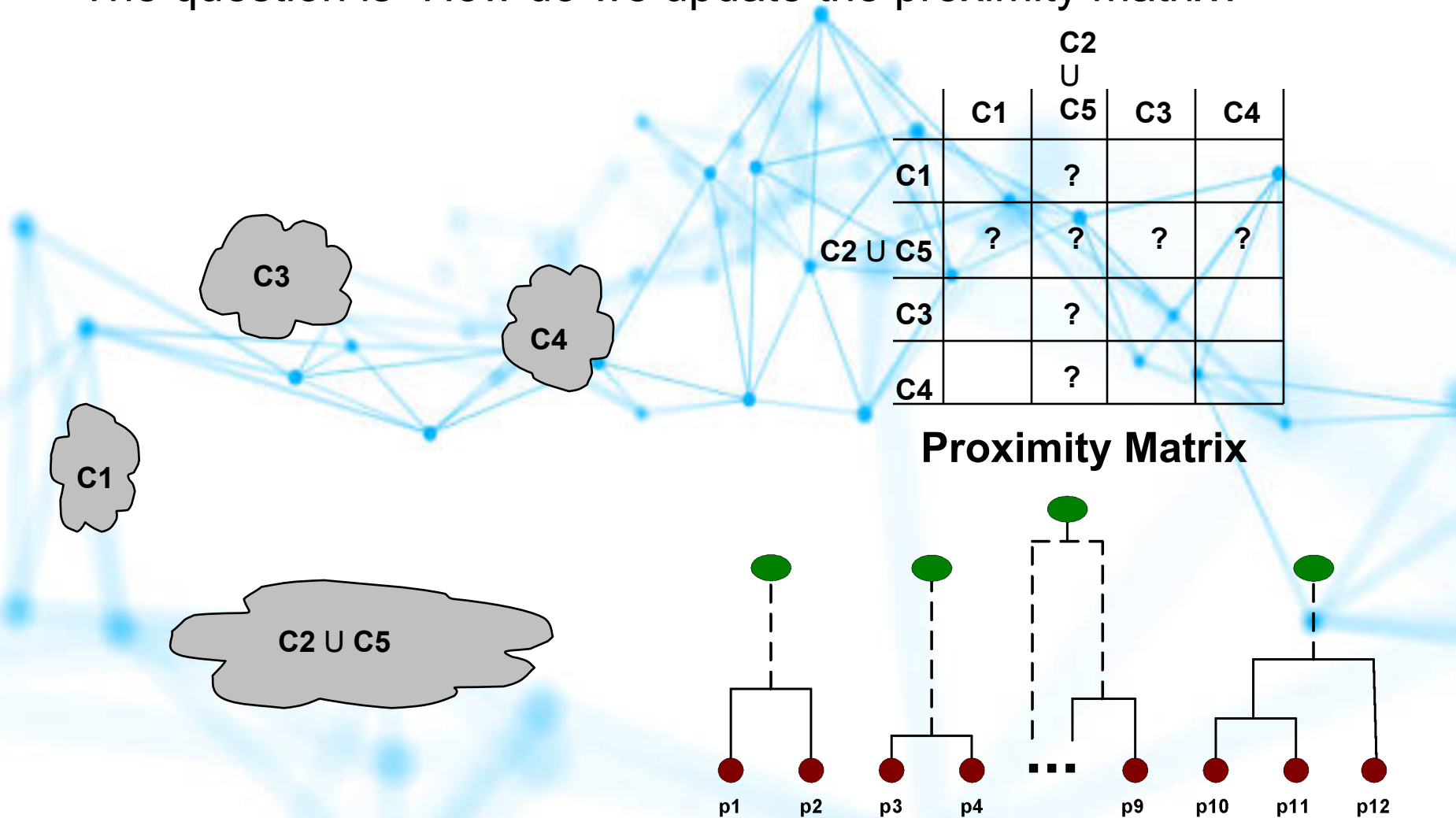
	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Proximity Matrix

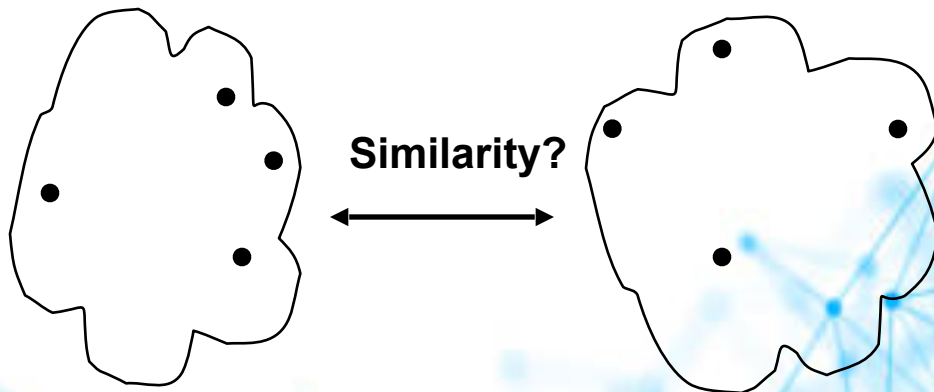


After Merging

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



How to Define Inter-Cluster Similarity

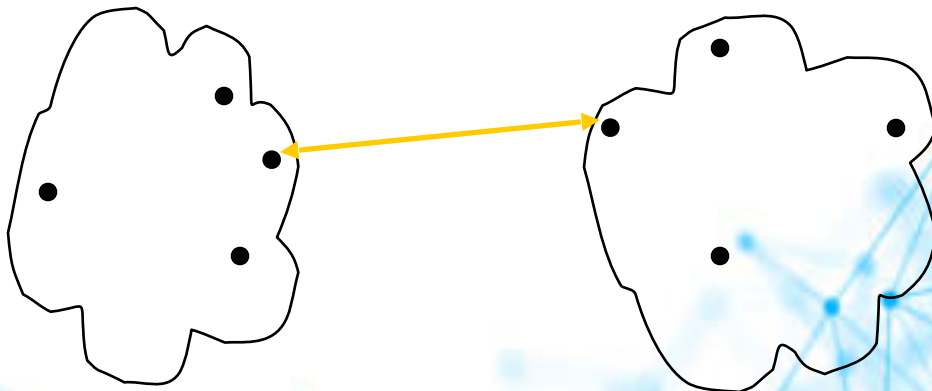


	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						
.						

Proximity Matrix

- MIN
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
 - Ward's Method uses squared error

How to Define Inter-Cluster Similarity

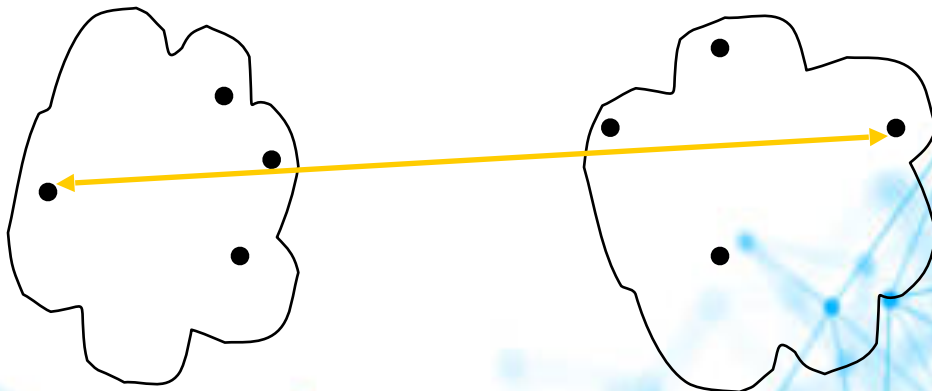


	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						
.						

- **MIN (Single Linkage)**
- **MAX (Complete Linkage)**
- **Group Average (Average Linkage)**
- **Distance Between Centroids**
- **Other methods driven by an objective function**
 - Ward's Method uses squared error

Proximity Matrix

How to Define Inter-Cluster Similarity

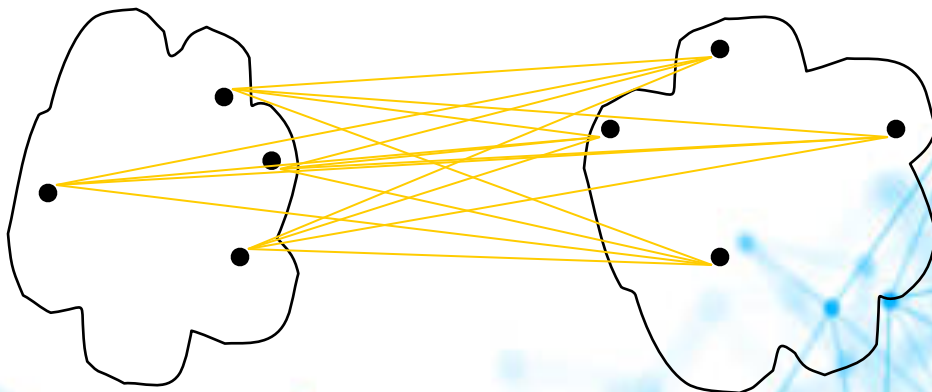


- MIN (Single Linkage)
- **MAX (Complete Linkage)**
- Group Average (Average Linkage)
- Distance Between Centroids
- Other methods driven by an objective function
 - Ward's Method uses squared error

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

Proximity Matrix

How to Define Inter-Cluster Similarity

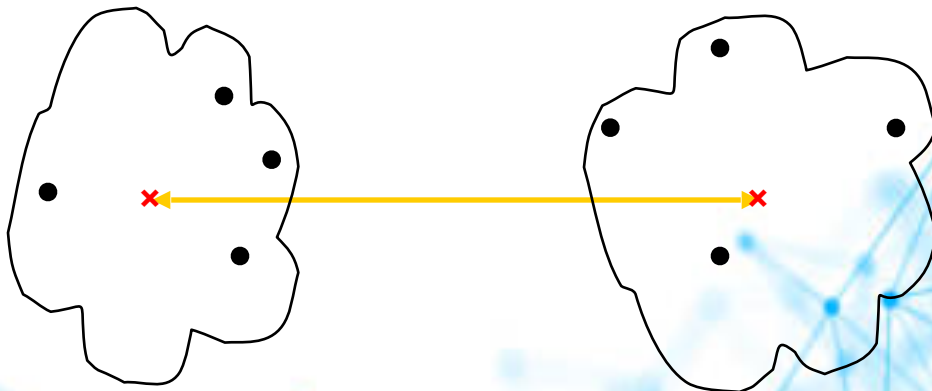


- MIN (Single Linkage)
- MAX (Complete Linkage)
- **Group Average (Average Linkage)**
- Distance Between Centroids
- Other methods driven by an objective function
 - Ward's Method uses squared error

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

Proximity Matrix

How to Define Inter-Cluster Similarity



- MIN (Single Linkage)
- MAX (Complete Linkage)
- Group Average (Average Linkage)
- **Distance Between Centroids**
- Other methods driven by an objective function
 - Ward's Method uses squared error

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

Proximity Matrix

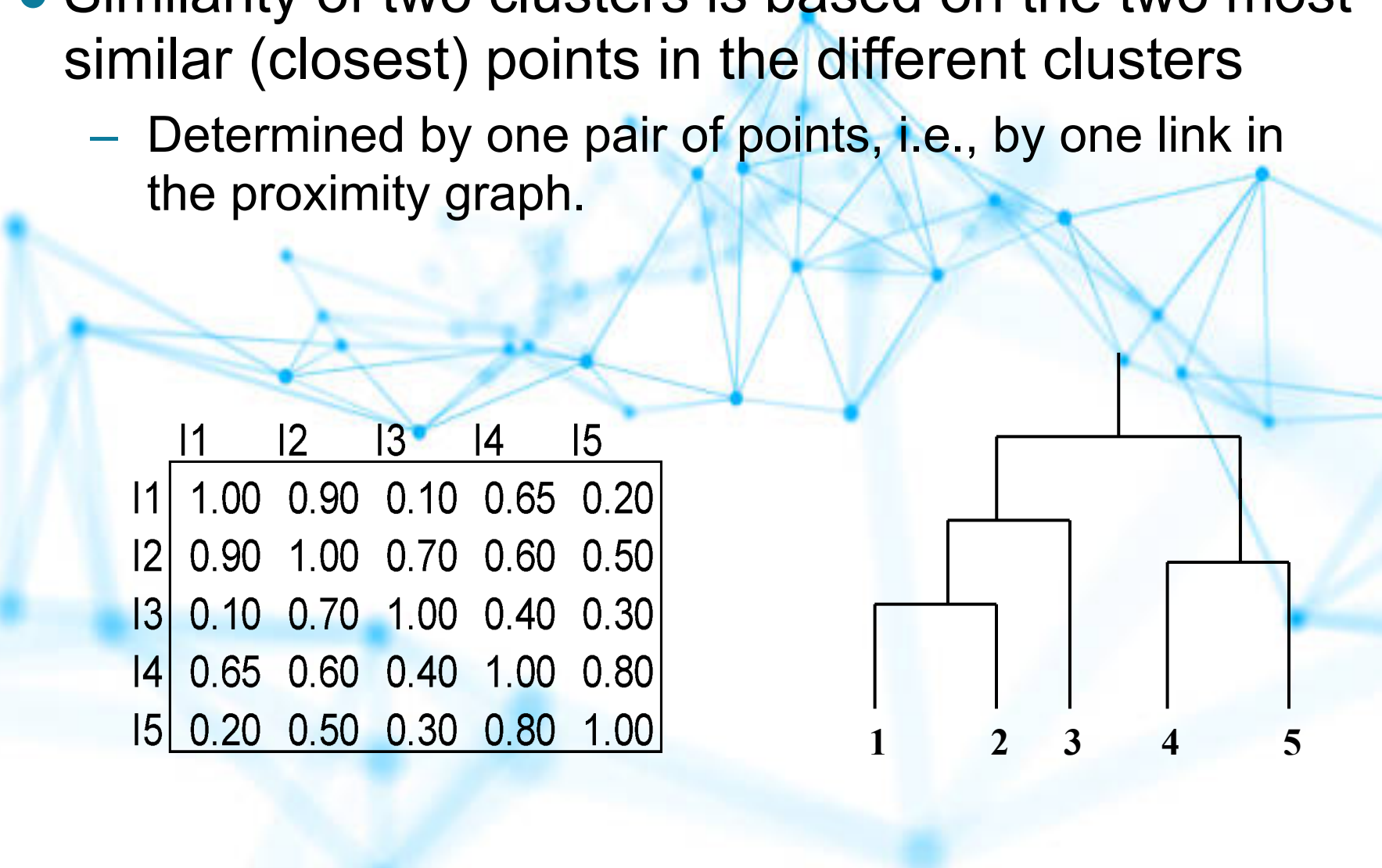
How to Define Inter-Cluster Similarity



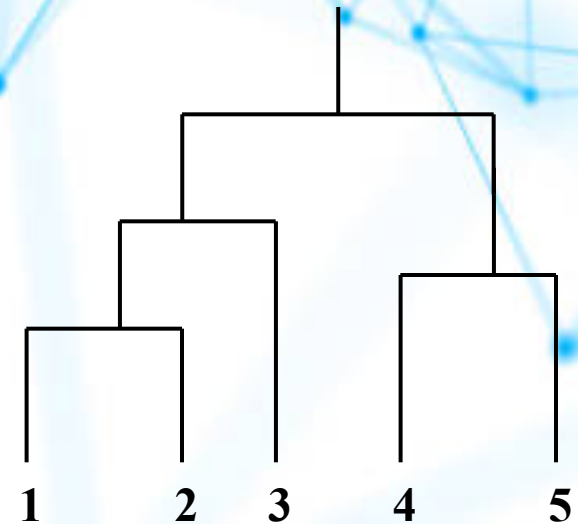
Single Linkage	This is the distance between the closest members of the two clusters.
Complete Linkage	This is the distance between the members that are farthest apart.
Average Linkage	This method involves looking at the distances between all pairs and averages all of these distances. This is also called Unweighted Pair Group Mean Averaging.

Cluster Similarity: MIN or Single Link

- Similarity of two clusters is based on the two most similar (closest) points in the different clusters
 - Determined by one pair of points, i.e., by one link in the proximity graph.



	I1	I2	I3	I4	I5
I1	1.00	0.90	0.10	0.65	0.20
I2	0.90	1.00	0.70	0.60	0.50
I3	0.10	0.70	1.00	0.40	0.30
I4	0.65	0.60	0.40	1.00	0.80
I5	0.20	0.50	0.30	0.80	1.00



Cluster Similarity: MIN or Single Link

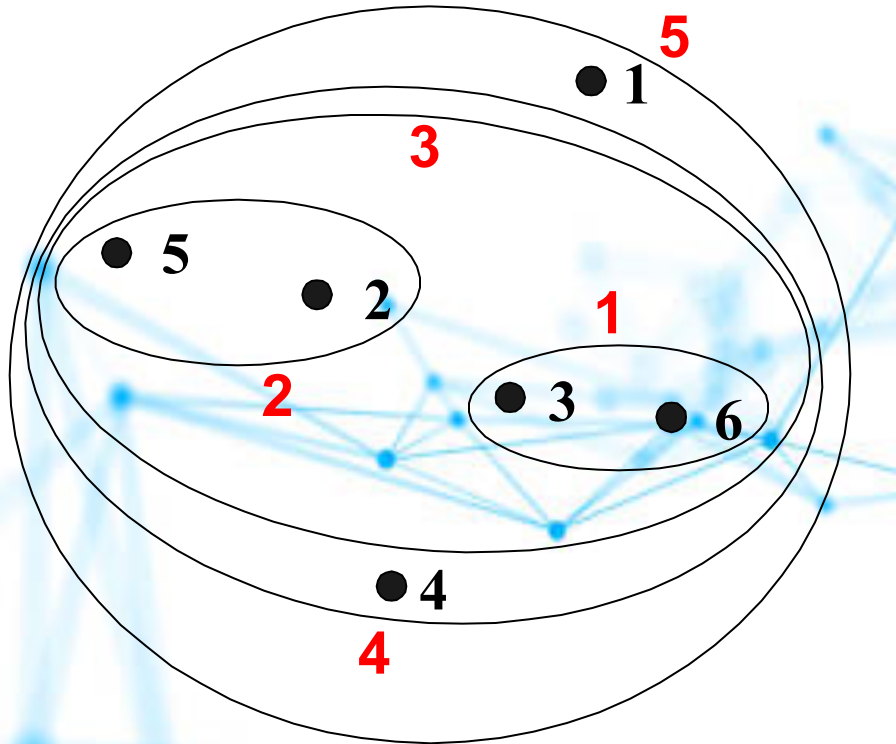
	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

Euclidean distance matrix for 6 points.

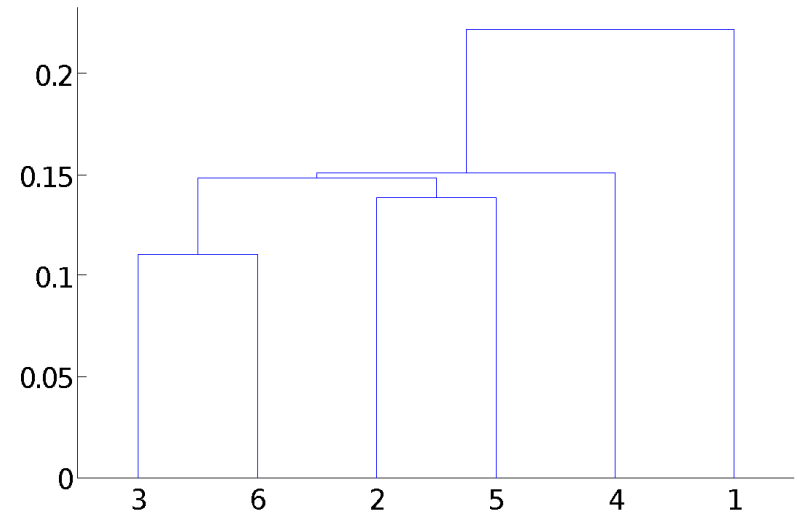
Hierarchical Clustering: MIN

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

Table 8.4. Euclidean distance matrix for 6 points.



Nested Clusters



Dendrogram

Cluster Similarity: MIN or Single Link

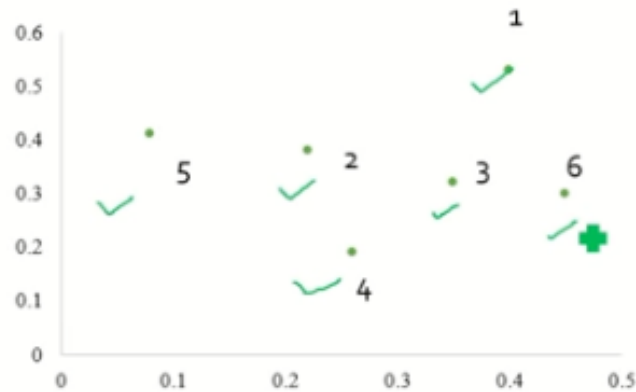


$$\begin{aligned} \text{dist}(\{3, 6\}, \{2, 5\}) &= \min(\text{dist}(3, 2), \text{dist}(6, 2), \text{dist}(3, 5), \text{dist}(6, 5)) \\ &= \min(0.15, 0.25, 0.28, 0.39) \\ &= 0.15. \end{aligned}$$

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

Table 8.4. Euclidean distance matrix for 6 points.

	X	Y
P1	0.40	0.53
P2	0.22	0.38
P3	0.35	0.32
P4	0.26	0.19
P5	0.08	0.41
P6	0.45	0.30



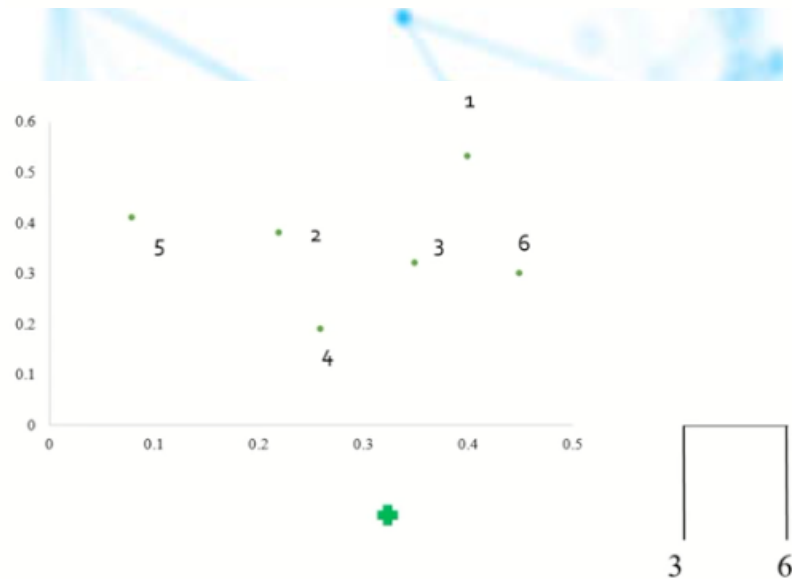
- Calculate Euclidean distance, create the distance matrix.

$$\text{Distance} [(x,y), (a,b)] = \sqrt{(x-a)^2 + (y-b)^2}$$

$$\begin{aligned} \text{Distance (P1,P2)} &= \sqrt{(0.40 - 0.22)^2 + (0.53 - 0.38)^2} \\ (0.40, 0.53), (0.22, 0.38) &= \sqrt{(0.18)^2 + (0.15)^2} \\ &= \sqrt{0.0324 + 0.0225} \\ &= \sqrt{0.0549} \\ &= 0.23 \end{aligned}$$

- The distance matrix is

	P1	P2	P3	P4	P5	P6
P1	0					
P2	0.23	0				
P3	0.22	0.15	0			
P4	0.37	0.20	0.15	0		
P5	0.34	0.14	0.28	0.29	0	
P6	0.23	0.25	0.11	0.22	0.39	0



- The distance matrix is

	P1	P2	P3	P4	P5	P6
P1	0					
P2	0.23	0				
P3	0.22	0.15	0			
P4	0.37	0.20	0.15	0		
P5	0.34	0.14	0.28	0.29	0	
P6	0.23	0.25	0.11	0.22	0.39	0

- To update the distance matrix $\text{MIN}[\text{dist}(P3, P6), P1]$

- $\text{MIN}(\text{dist}(P3, P1), (P6, P1))$

$$= \min[(0.22, 0.23)]$$

$$= 0.22$$

- To update the distance matrix $\text{MIN}[\text{dist}(P3, P6), P2]$

- $\text{MIN}(\text{dist}(P3, P2), (P6, P2))$

$$= \min[(0.15, 0.25)]$$

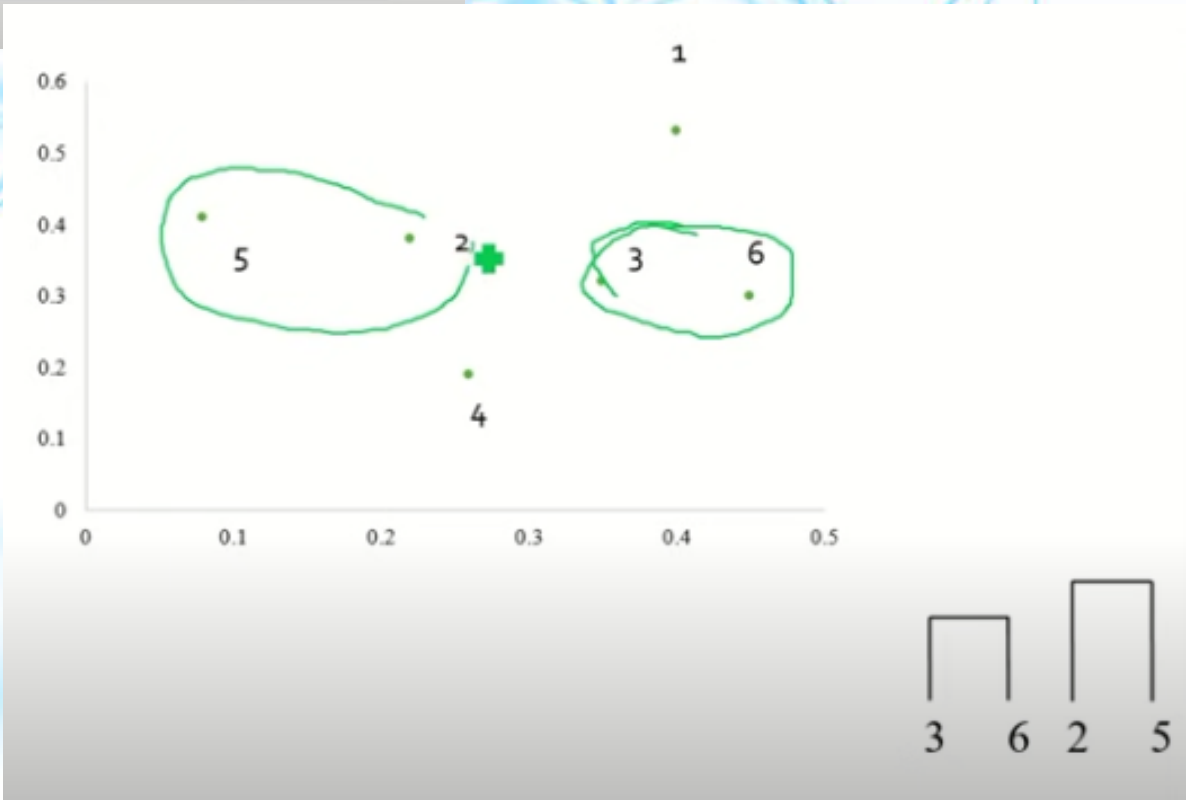
$$= 0.15$$

- The updated distance matrix for cluster P3, P6

	P1	P2	P3,P6	P4	P5
P1	0				
P2	0.23	0			
P3,P6	0.22	0.15	0		
P4	0.37	0.20	0.15	0	
P5	0.34	0.14	0.28	0.29	0

■ The distance matrix is

	P1	P2	P3,P6	P4	P5
P1	0				
P2	0.23	0			
P3,P6	0.22	0.15	0		
P4	0.37	0.20	0.15	0	
P5	0.34	0.14	0.28	0.29	0



- To update the distance matrix $\text{MIN}[\text{dist}(\text{P2}, \text{P5}), \text{P1}]$

- $\text{MIN}[\text{dist}(\text{P2}, \text{P1}), (\text{P5}, \text{P1})]$

$$= \min[(0.23, 0.34)]$$

$$= 0.23$$

- To update the distance matrix $\text{MIN}[\text{dist}(\text{P2}, \text{P5}), (\text{P3}, \text{P6})]$

- $\text{MIN}[\text{dist}(\text{P2}, (\text{P3}, \text{P6})), (\text{P5}, (\text{P3}, \text{P6}))]$

$$= \min[(0.15, 0.28)]$$

$$= 0.15$$

- To update the distance matrix $\text{MIN}[\text{dist}(\text{P2}, \text{P5}), \text{P4}]$

- $\text{MIN}[\text{dist}(\text{P2}, \text{P4}), (\text{P5}, \text{P4})]$

$$= \min[(0.20, 0.29)]$$

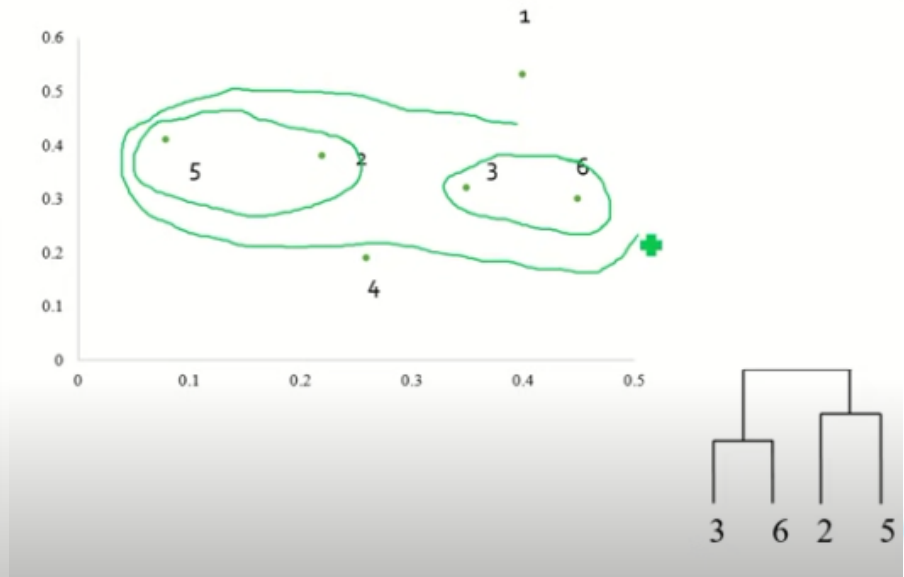
$$= 0.20$$

- The updated distance matrix for cluster P2,P5

	P1	P2,P5	P3,P6	P4
P1	0			
P2,P5	0.23	0		
P3,P6	0.22	0.15	0	
P4	0.37	0.20	0.15	0

- The distance matrix is

	P1	P2,P5	P3,P6	P4
P1	0			
P2,P5	0.23	0		
P3,P6	0.22	0.15	0	
P4	0.37	0.20	0.15	0



- To update the distance matrix $\text{MIN}[\text{dist}((P2,P5),(P3,P6)),P1]$

- $\text{MIN}[\text{dist}((P2,P5),P1), ((P3,P6),P1)]$
 $= \min[(0.23,0.22)]$
 $= 0.22$

- To update the distance matrix $\text{MIN}[\text{dist}((P2,P5),(P3,P6)),P4]$

- $\text{MIN}[\text{dist}((P2,P5),P4), ((P3,P6),P4)]$
 $= \min[(0.20,0.15)]$
 $= 0.15$

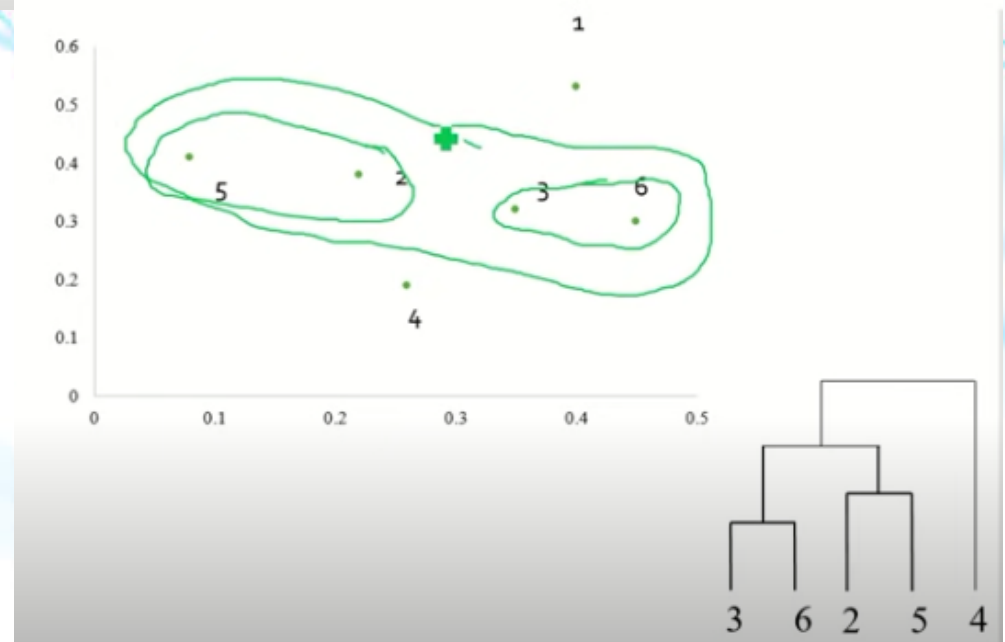
- The updated distance matrix for cluster P2,P5,P3,P6

	P1	P2,P5,P3,P6	P4
P1	0		
P2,P5,P3,P6	0.22	0	
P4	0.37	0.15	0



- The distance matrix is

	P1	P2,P5,P3,P6	P4
P1	0		
P2,P5,P3,P6	0.22	0	
P4	0.37	0.15	0



- To update the distance matrix $\text{MIN}[\text{dist}(\text{P2}, \text{P5}, \text{P3}, \text{P6}), \text{P4}]$

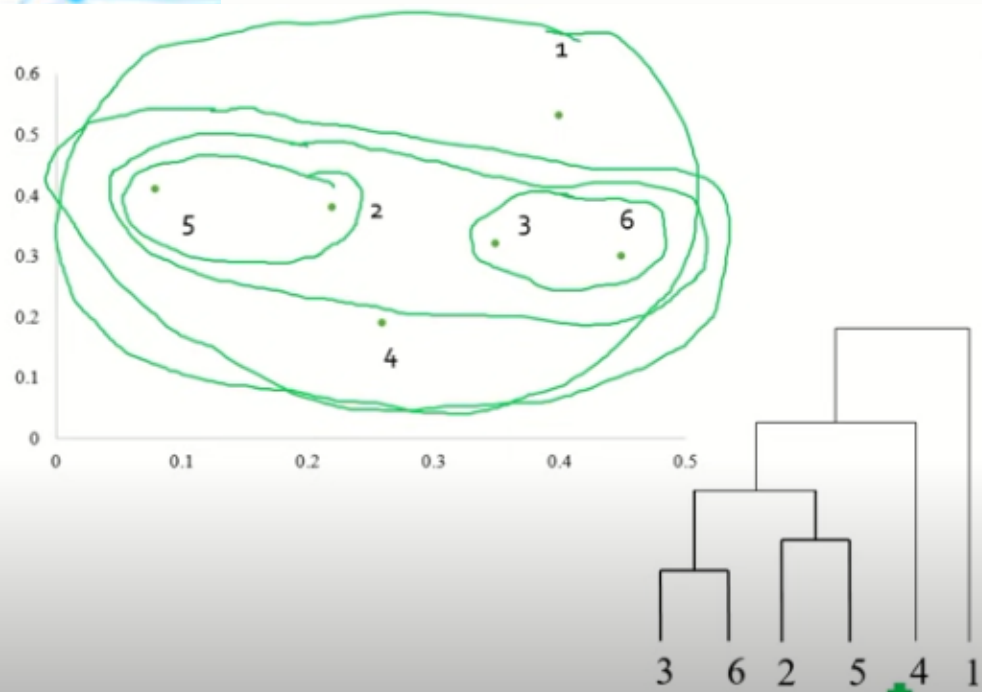
- $\text{MIN}[\text{dist}((\text{P2}, \text{P5}, \text{P3}, \text{P6}), \text{P1}), (\text{P4}, \text{P1})]$

$= \min[(0.22, 0.37)]$

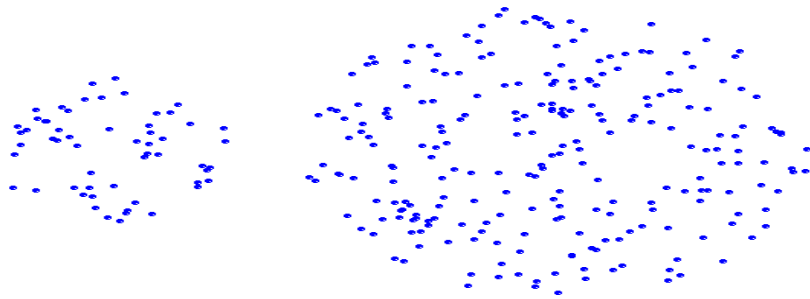
$= 0.22$

- The updated distance matrix for cluster $\text{P2}, \text{P5}, \text{P3}, \text{P6}, \text{P4}$

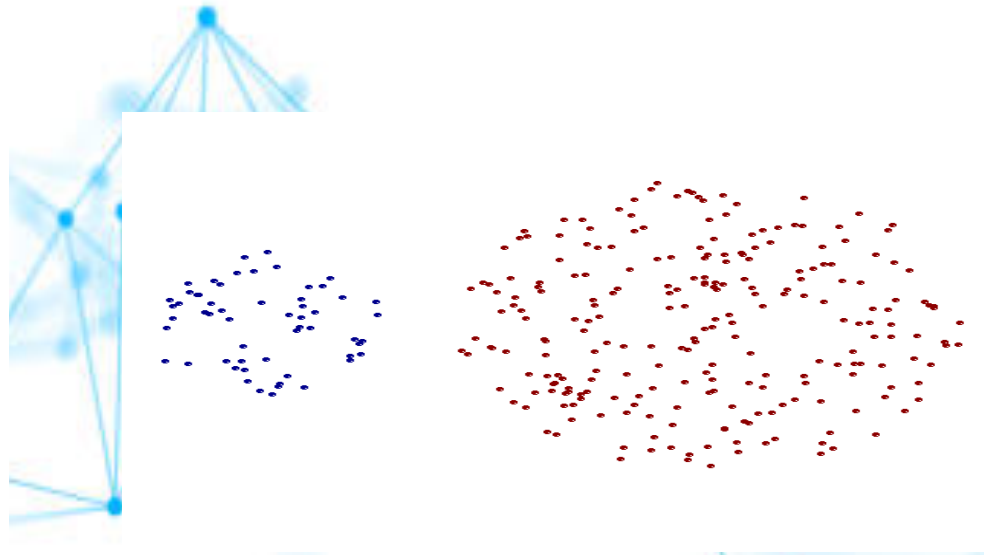
	P1	P2, P5, P3, P6, P4
P1	0	
P2, P5, P3, P6, P4	0.22	0



Strength of MIN



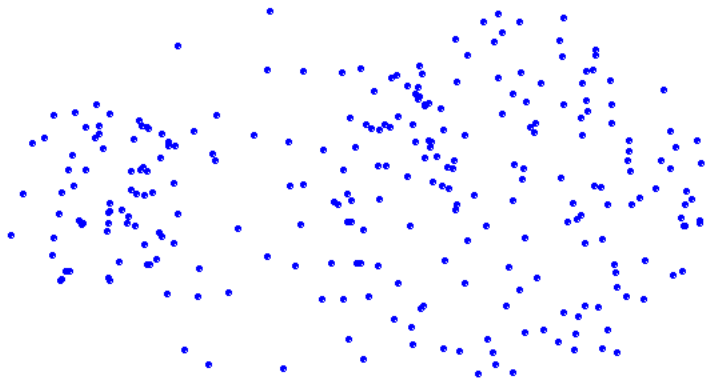
Original Points



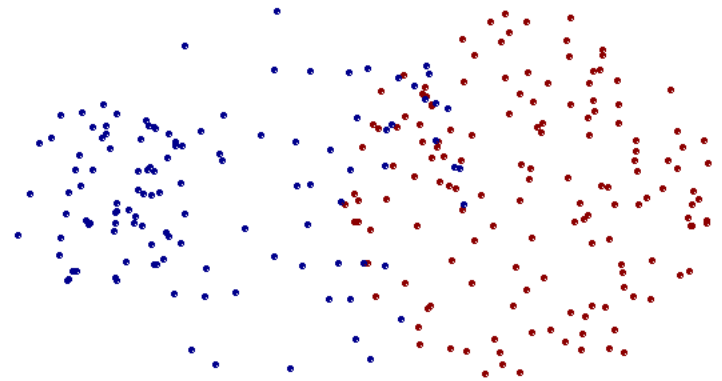
Two Clusters

- Can handle non-elliptical shapes

Limitations of MIN



Original Points

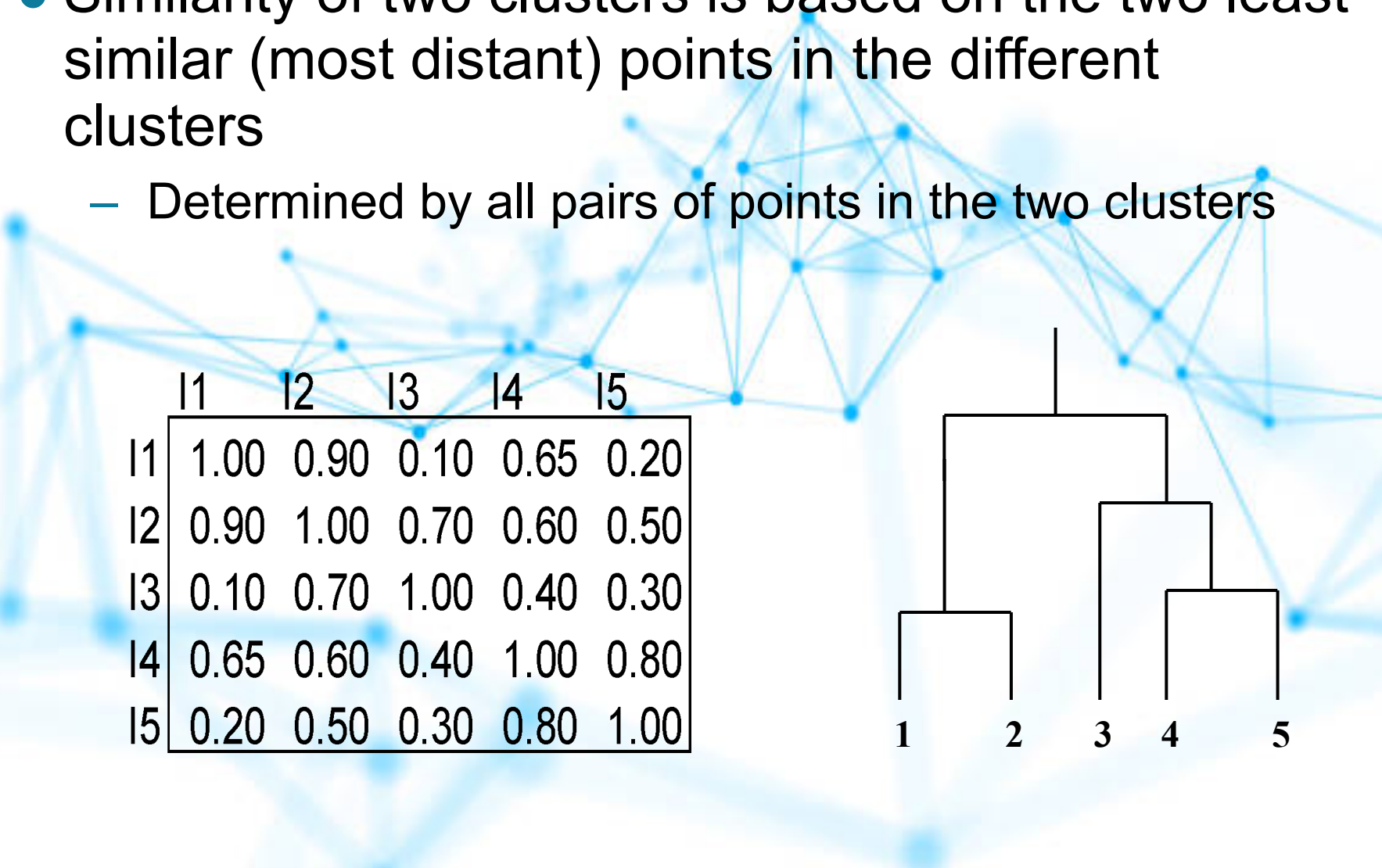


Two Clusters

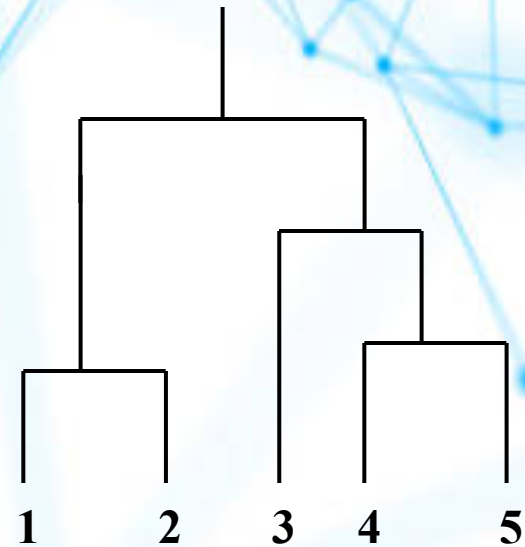
- **Sensitive to noise and outliers**

Cluster Similarity: MAX or Complete Linkage

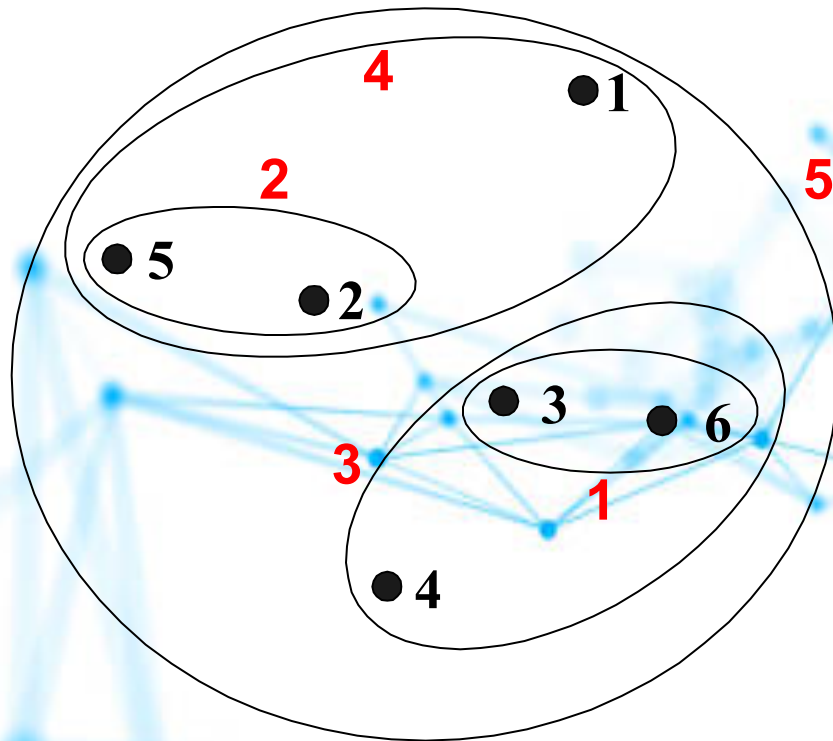
- Similarity of two clusters is based on the two least similar (most distant) points in the different clusters
 - Determined by all pairs of points in the two clusters



	I1	I2	I3	I4	I5
I1	1.00	0.90	0.10	0.65	0.20
I2	0.90	1.00	0.70	0.60	0.50
I3	0.10	0.70	1.00	0.40	0.30
I4	0.65	0.60	0.40	1.00	0.80
I5	0.20	0.50	0.30	0.80	1.00



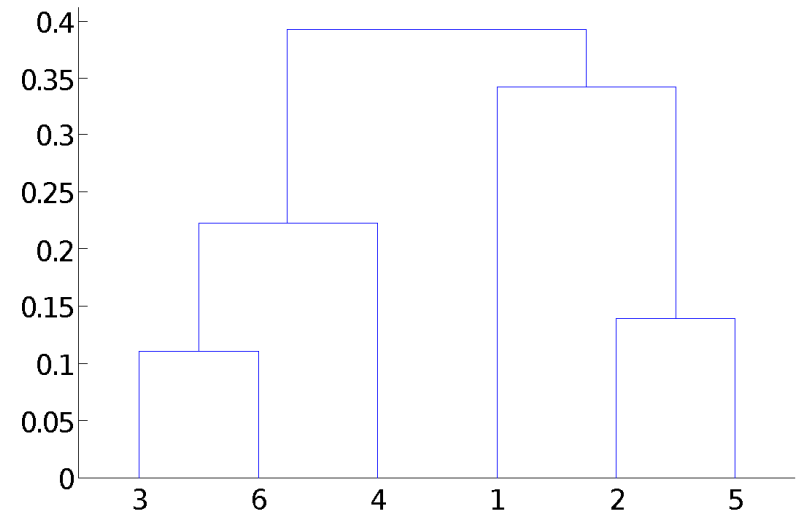
Hierarchical Clustering: MAX



Nested Clusters

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

Table 8.4. Euclidean distance matrix for 6 points.



Dendrogram

Cluster Similarity: MAX

$$\begin{aligned} \text{dist}(\{3, 6\}, \{4\}) &= \max(\text{dist}(3, 4), \text{dist}(6, 4)) \\ &= \max(0.15, 0.22) \\ &= 0.22. \end{aligned}$$

$$\text{dist}(\{3, 6\}, \{2, 5\})$$

$$\text{dist}(\{3, 6\}, \{1\})$$

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

Table 8.4. Euclidean distance matrix for 6 points.

- To update the distance matrix $\text{MAX}[\text{dist}(P3,P6),P1)]$



- $\text{MAX}(\text{dist}(P3,P1), (P6,P1))$

$$= \text{MAX}[(0.22,0.23)]$$

$$= 0.23$$

- To update the distance matrix $\text{MAX}[\text{dist}(P3,P6),P2)]$

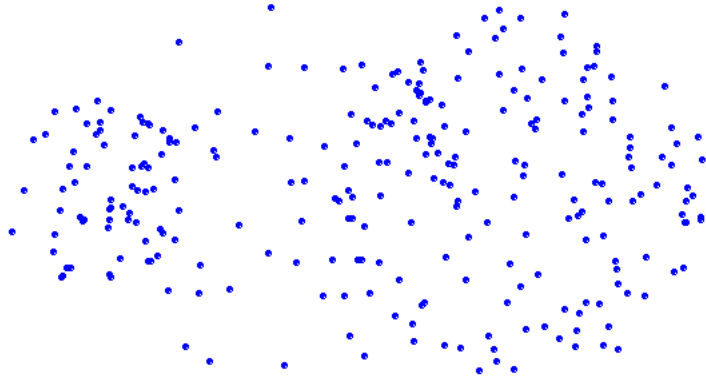
- $\text{MAX}(\text{dist}(P3,P2), (P6,P2))$

$$= \text{MAX}[(0.15,0.25)]$$

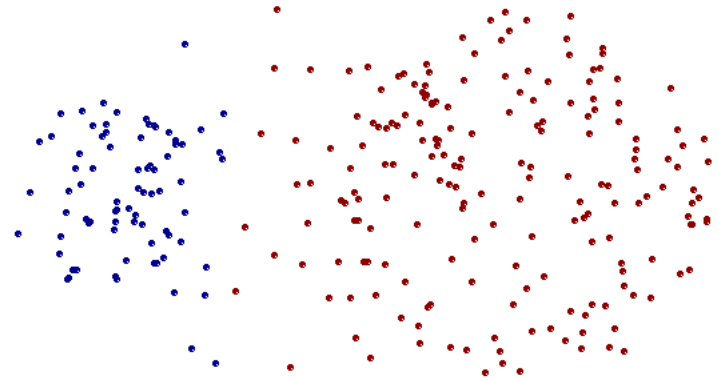
$$= 0.25$$

	P1	P2	P3	P4	P5	P6
P1	0					
P2	0.23	0				
P3	0.22	0.15	0			
P4	0.37	0.20	0.15	0		
P5	0.34	0.14	0.28	0.29	0	
P6	0.23	0.25	0.11	0.22	0.39	0

Strength of MAX



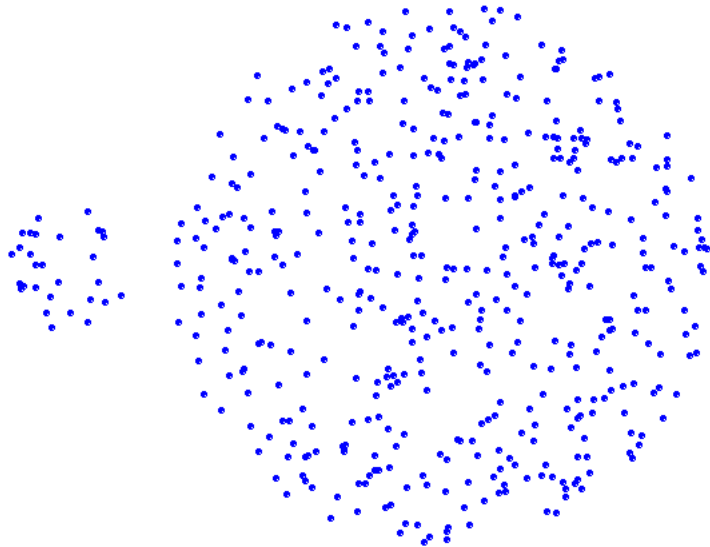
Original Points



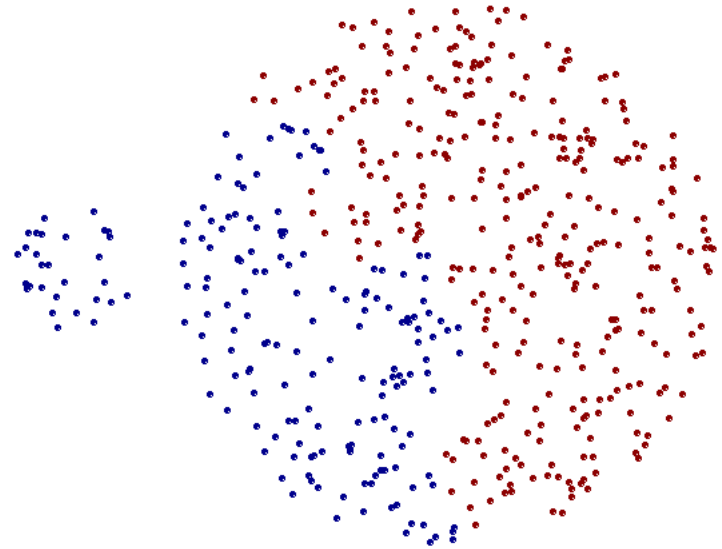
Two Clusters

- **Less susceptible to noise and outliers**

Limitations of MAX



Original Points



Two Clusters

- Tends to break large clusters
- Biased towards globular clusters

Cluster Similarity: Group Average

- Proximity of two clusters is the average of pairwise proximity between points in the two clusters.

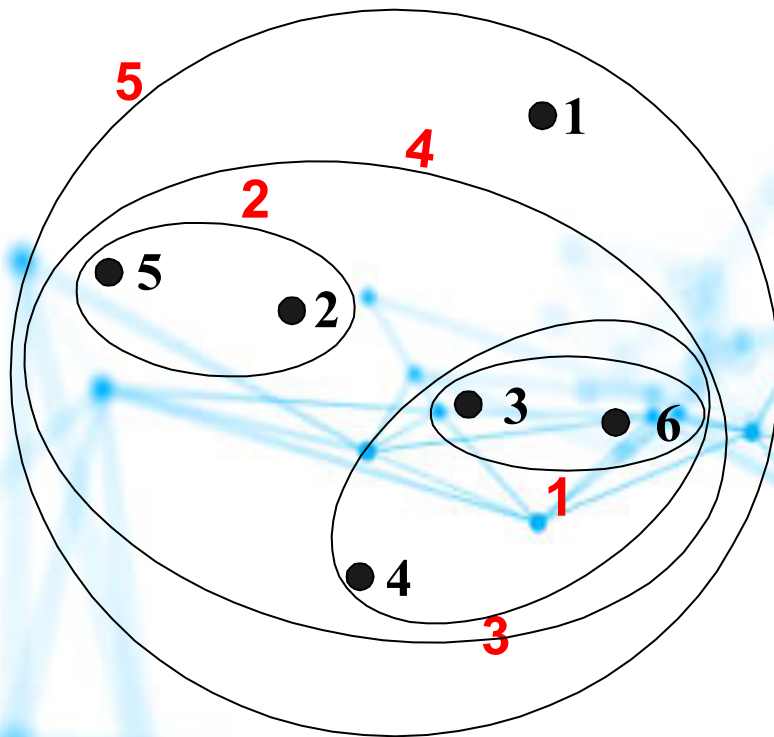
$$\text{proximity}(\text{Cluster}_i, \text{Cluster}_j) = \frac{\sum_{\substack{p_i \in \text{Cluster}_i \\ p_j \in \text{Cluster}_j}} \text{proximity}(p_i, p_j)}{|\text{Cluster}_i| * |\text{Cluster}_j|}$$

- Need to use average connectivity for scalability since total proximity favors large clusters

	I1	I2	I3	I4	I5
I1	1.00	0.90	0.10	0.65	0.20
I2	0.90	1.00	0.70	0.60	0.50
I3	0.10	0.70	1.00	0.40	0.30
I4	0.65	0.60	0.40	1.00	0.80
I5	0.20	0.50	0.30	0.80	1.00



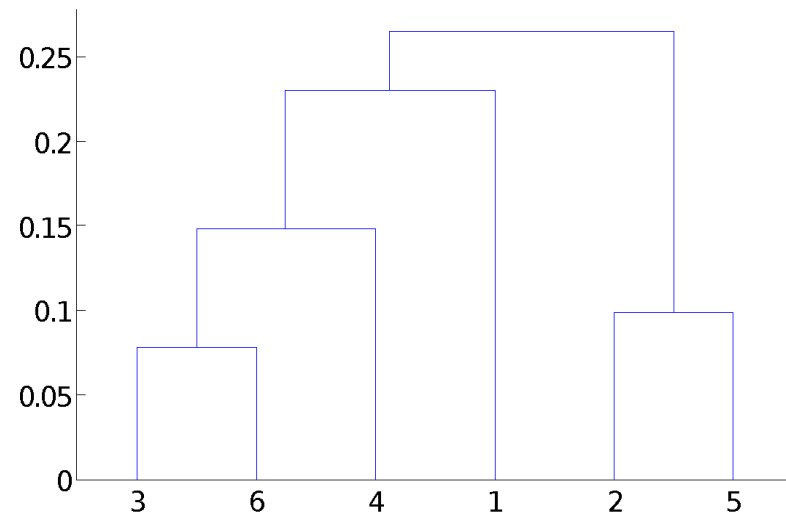
Hierarchical Clustering: Group Average



Nested Clusters

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

Table 8.4. Euclidean distance matrix for 6 points.

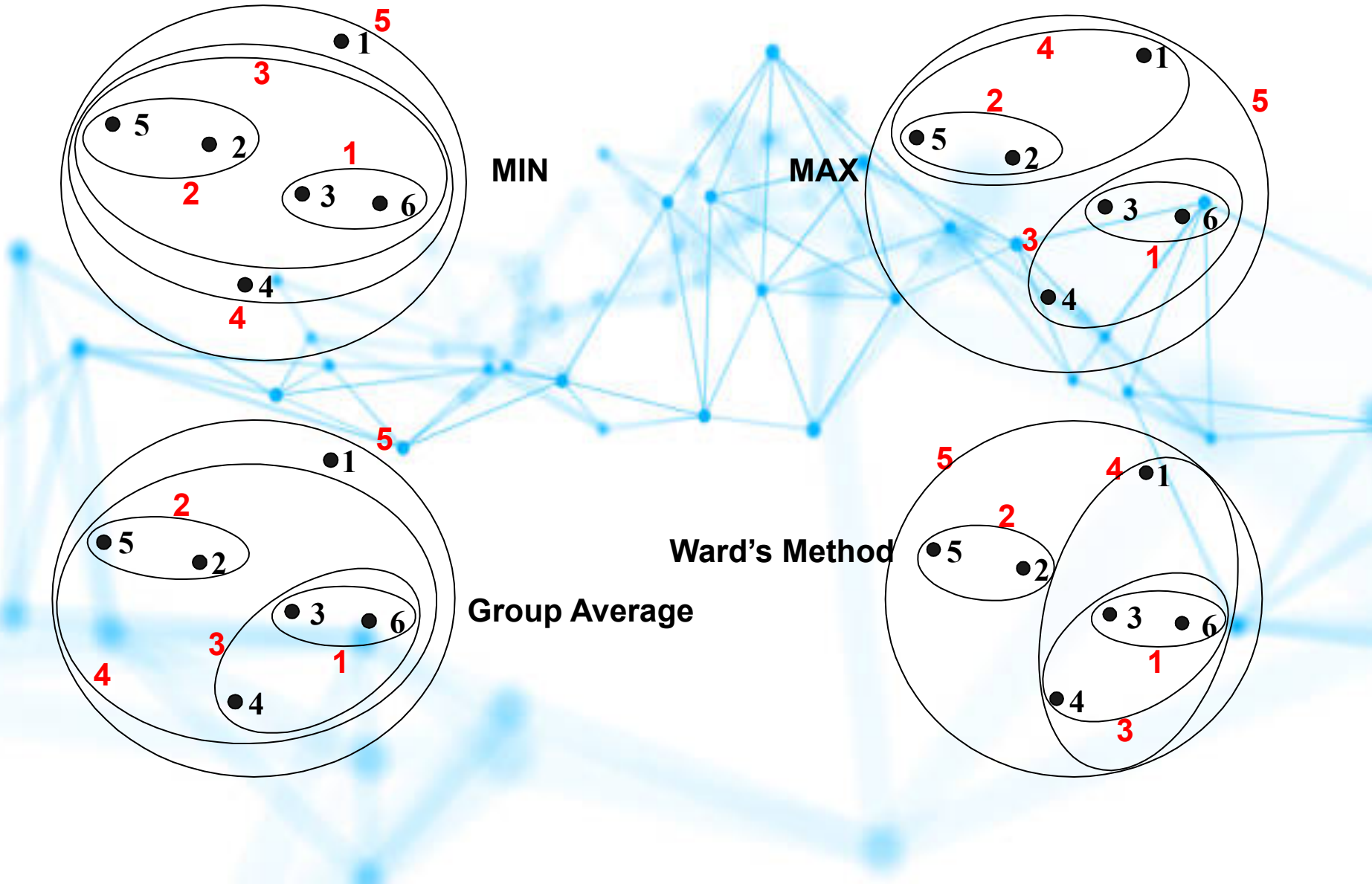


Dendrogram

Hierarchical Clustering: Group Average

- Compromise between Single and Complete Link
- Strengths
 - Less susceptible to noise and outliers
- Limitations
 - Biased towards globular clusters

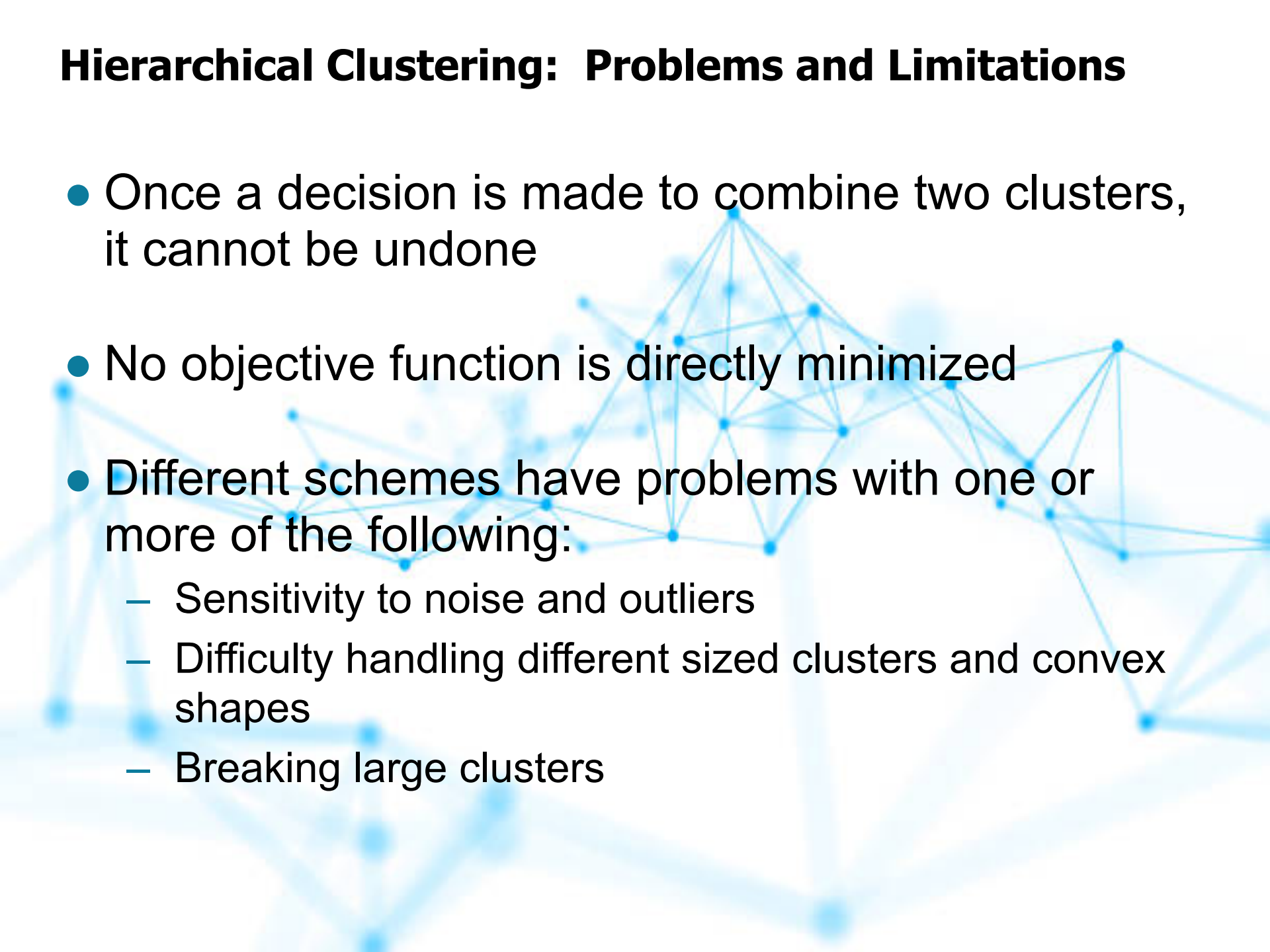
Hierarchical Clustering: Comparison



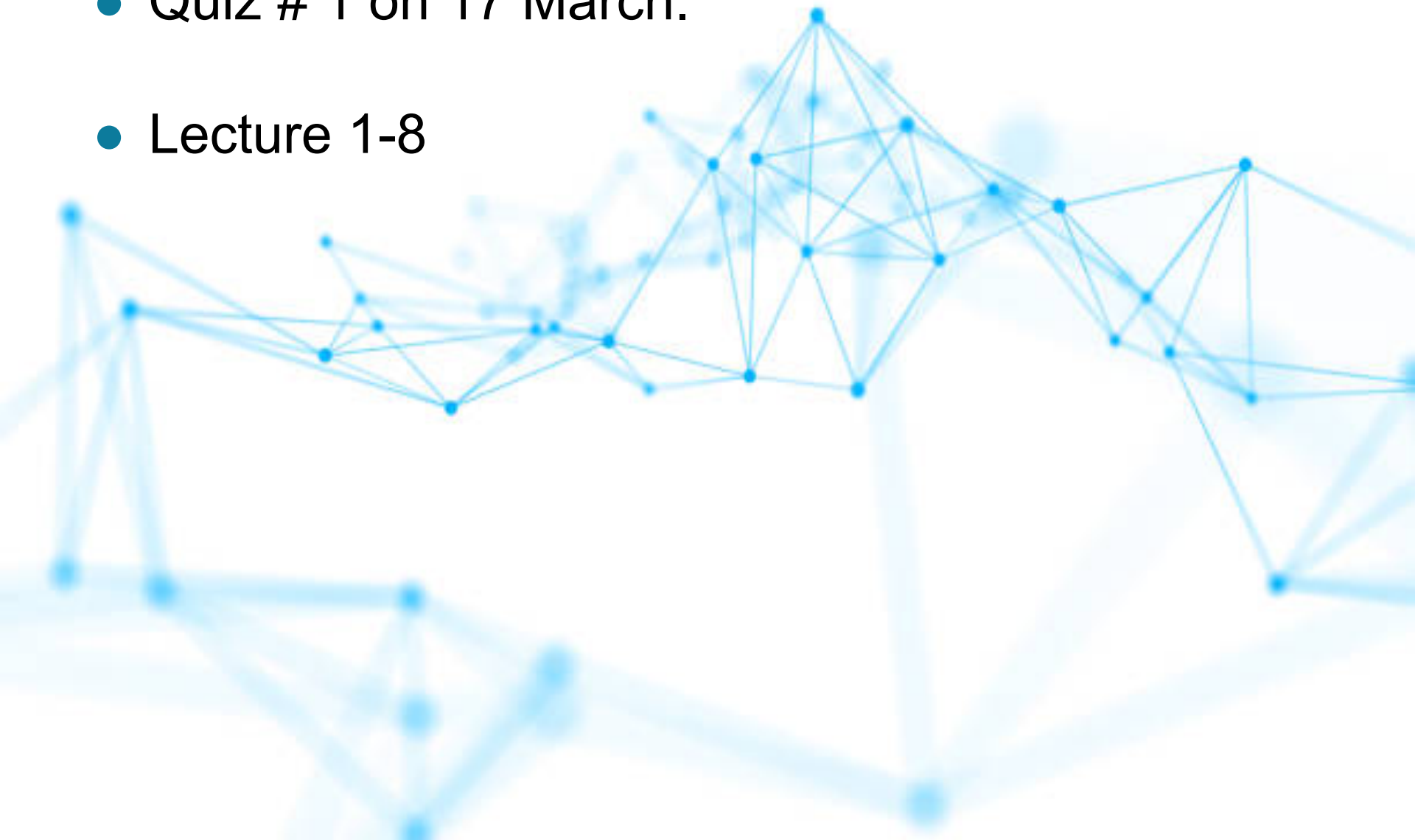
Hierarchical Clustering: Time and Space requirements

- $O(N^2)$ space since it uses the proximity matrix.
 - N is the number of 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 for some approaches

Hierarchical Clustering: Problems and Limitations

- Once a decision is made to combine two clusters, it cannot be undone
 - No objective function is directly minimized
 - Different schemes have problems with one or more of the following:
 - Sensitivity to noise and outliers
 - Difficulty handling different sized clusters and convex shapes
 - Breaking large clusters
- 

- Quiz # 1 on 17 March.
- Lecture 1-8



Resources

- <https://www.youtube.com/watch?v=Cy3ci0Vqs3Y>
- <https://www.youtube.com/watch?v=RdT7bhm1M3E>
- <https://www.youtube.com/watch?v=9U4h6pZw6f8>