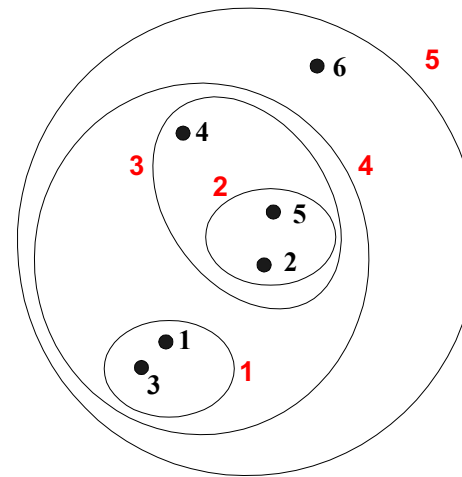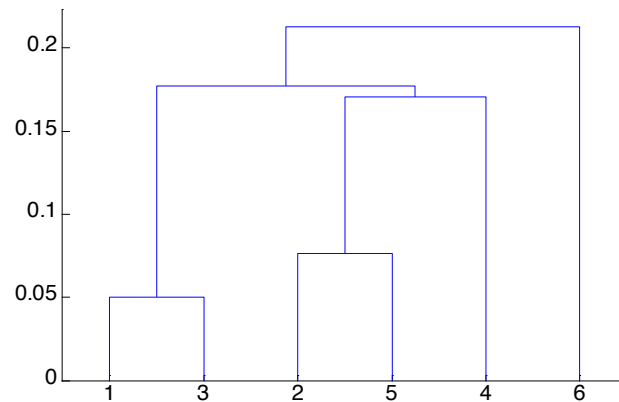# Hierarchical Clustering

# 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

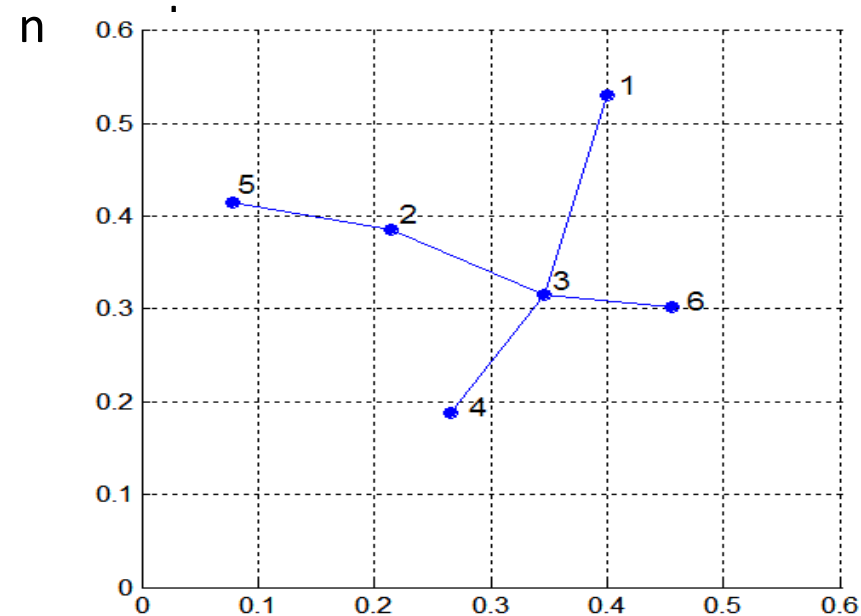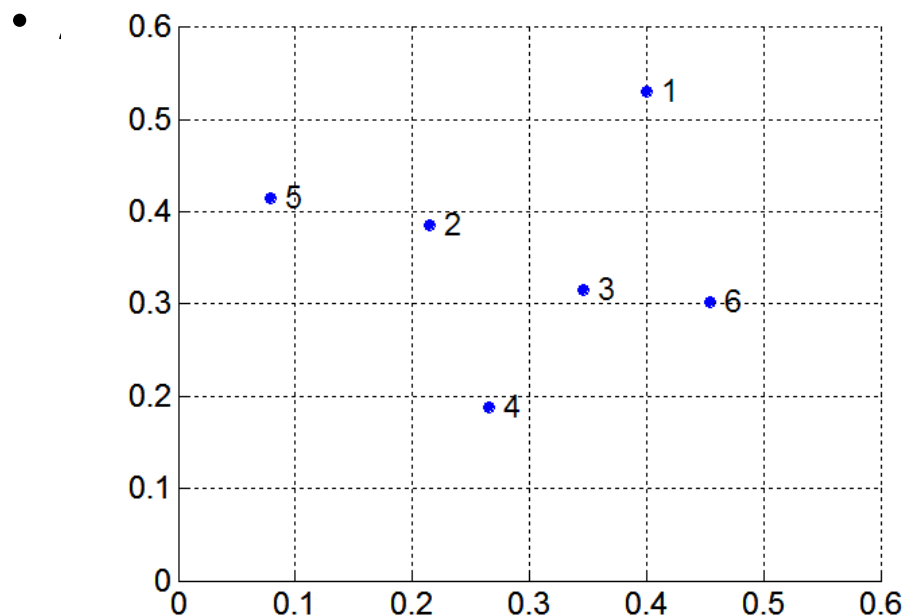# Strengths of Hierarchical Clustering

- Does 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

- The dendrogram may correspond to meaningful taxonomies

# 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

# MST: Divisive Hierarchical Clustering

- ## Build MST (Minimum Spanning Tree)
  - Start with a tree that consists of any point
  - In successive steps, look for the closest pair of points (p, q)  such that one point (p) is in the current tree but the other (q) is not

# MST: Divisive Hierarchical Clustering

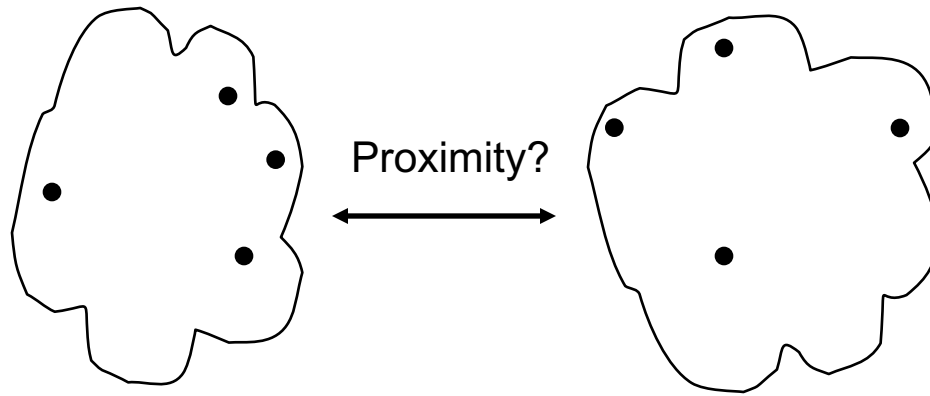- Use MST for constructing hierarchy of clusters

**Algorithm 7.5** MST Divisive Hierarchical Clustering Algorithm

1: Compute a minimum spanning tree for the proximity graph.
2: **repeat**
3:     Create a new cluster by breaking the link corresponding to the largest distance (smallest similarity).
4: **until** Only singleton clusters remain

# 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
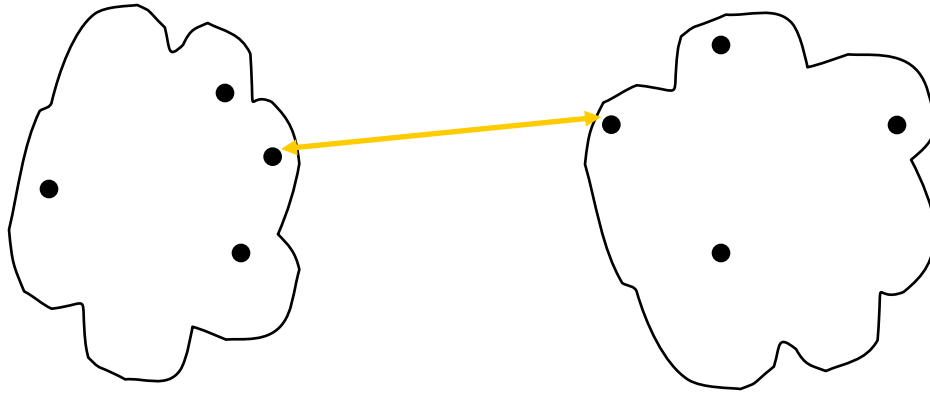
# How to Define Inter-Cluster Proximity



|     | p1 | p2 | p3 | p4 | p5 | . . . |
|-----|----|----|----|----|----|-------|
| p1  |    |    |    |    |    |       |
| p2  |    |    |    |    |    |       |
| p3  |    |    |    |    |    |       |
| p4  |    |    |    |    |    |       |
| p5  |    |    |    |    |    |       |
| .   |    |    |    |    |    |       |
| .   |    |    |    |    |    |       |
| .   |    |    |    |    |    |       |

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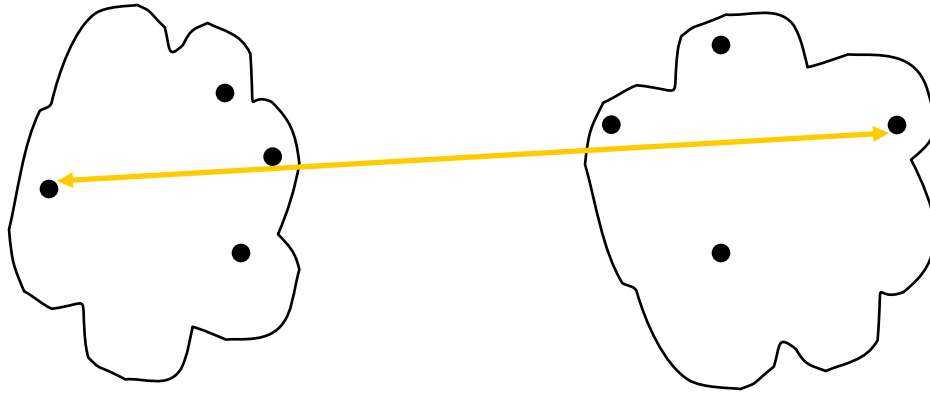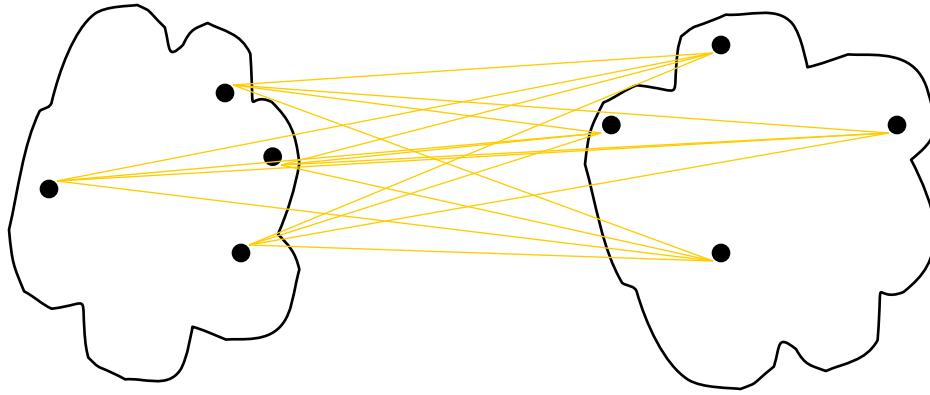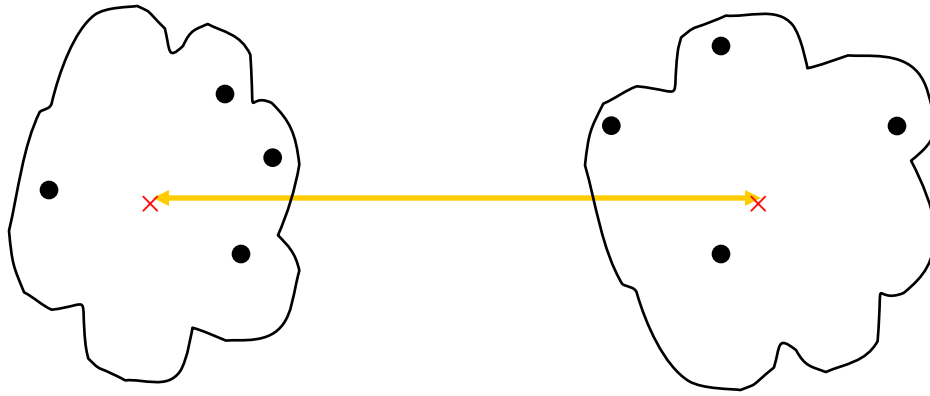
- MIN
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

# How to Define Inter-Cluster Proximity



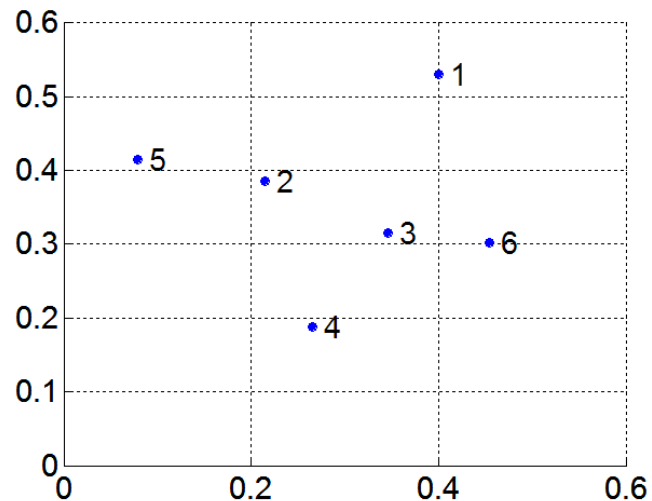|    | p1 | p2 | p3 | p4 | p5 | . . . |
|----|----|----|----|----|----|-------|
| p1 |    |    |    |    |    |       |
| p2 |    |    |    |    |    |       |
| p3 |    |    |    |    |    |       |
| p4 |    |    |    |    |    |       |
| p5 |    |    |    |    |    |       |
| .  |    |    |    |    |    |       |
| .  |    |    |    |    |    |       |
| .  |    |    |    |    |    |       |

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 Proximity

|     | 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 Proximity



| | 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 Proximity



|    | 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 changes in SSE

# MIN

- Proximity of two clusters is based on the two closest points in the different clusters
  - Determined by one pair of points

- Example:

Distance Matrix:

|     | 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 |

# Hierarchical Clustering: MIN



Nested Clusters

Dendrogram

# Strength of MIN

**Original Points**

**Six Clusters**

- **Can handle non-elliptical shapes**

# Limitations of MIN



**Two Clusters**

**Three Clusters**

**Original Points**

- **Sensitive to noise**

# MAX

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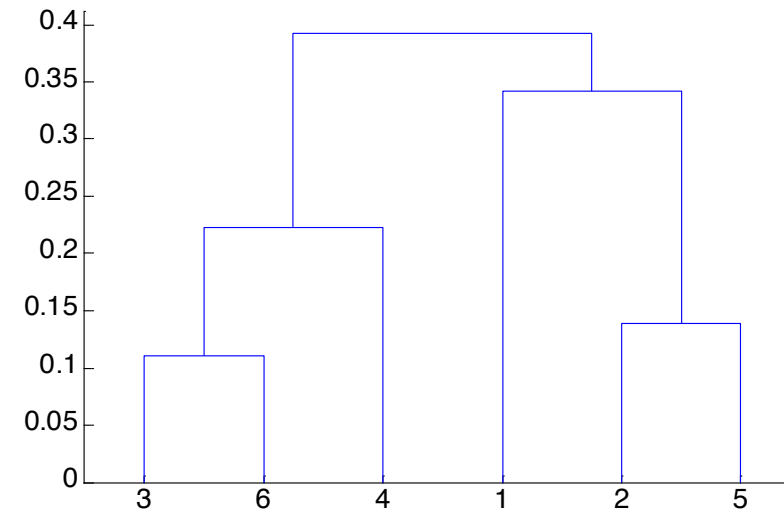  - Determined by all pairs of points in the two clusters

Distance Matrix:



|     | 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 |

# Hierarchical Clustering: MAX



Nested Clusters

Dendrogram

# Strength of MAX



**Original Points**　　　　　　　**Two Clusters**

- **Less susceptible to noise**

# Limitations of MAX



**Original Points**

**Two Clusters**

- **Tends to break large clusters**

# Group Average

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

$$\mathbf{proximity(Cluster_i, Cluster_j)} = \frac{\displaystyle\sum_{\substack{p_i \in Cluster_i \\ p_j \in Cluster_j}} \mathbf{proximity(p_i, p_j)}}{\mathbf{|Cluster_i| \times |Cluster_j|}}$$

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

Distance Matrix:

|     | 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 |

# Hierarchical Clustering: Group Average
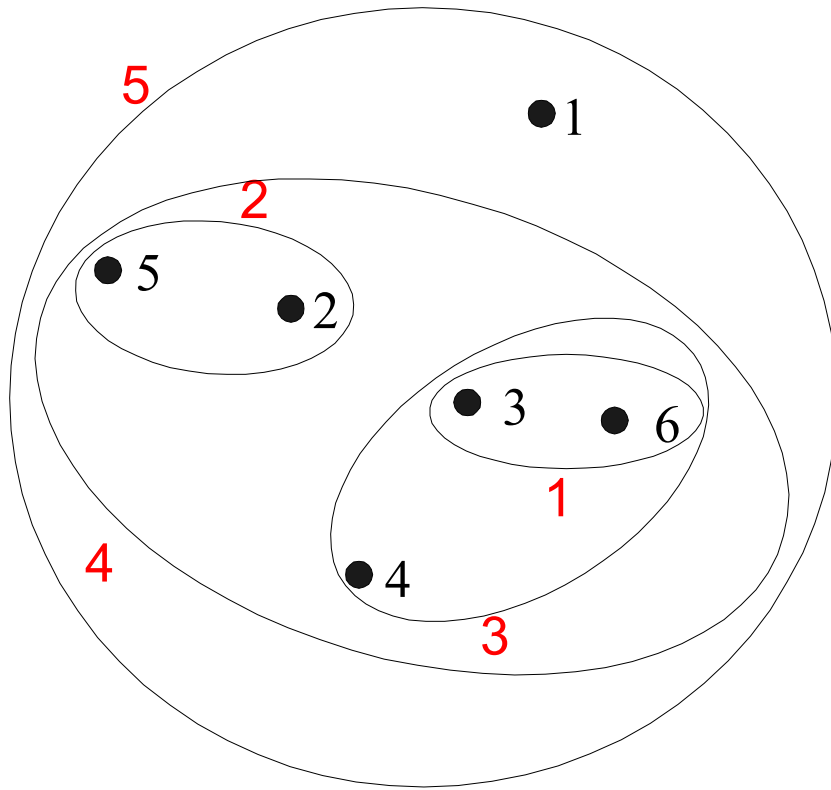


Nested Clusters

Dendrogram

# Hierarchical Clustering: Group Average

- Compromise between MIN and MAX

- Strengths
  - Less susceptible to noise

- Limitations
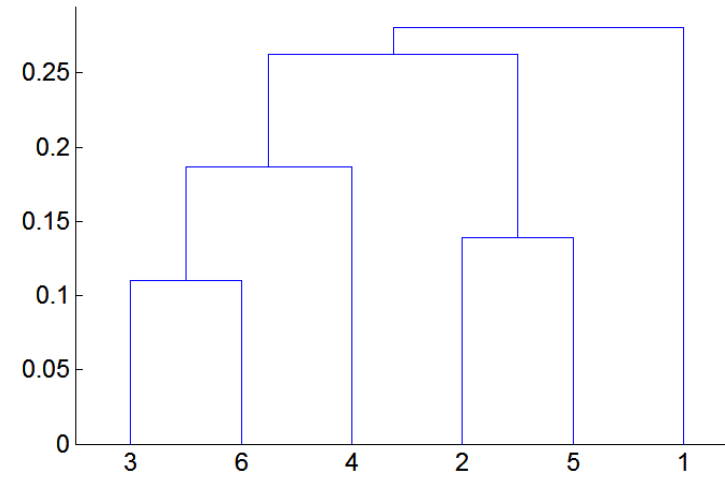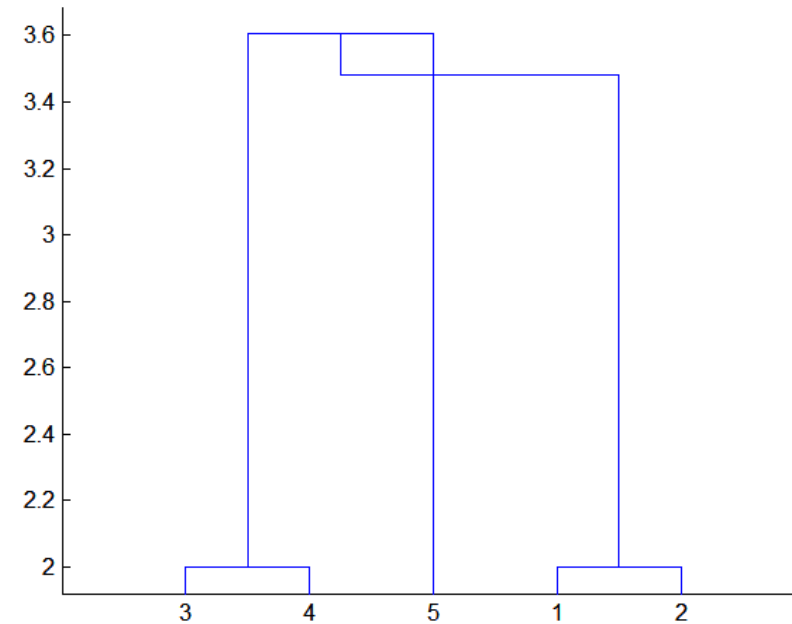  - Biased towards globular clusters

# Centroid Methods

- Similarity of two clusters is based on proximity between centroids of clusters

- Less susceptible to noise

- Biased towards globular clusters
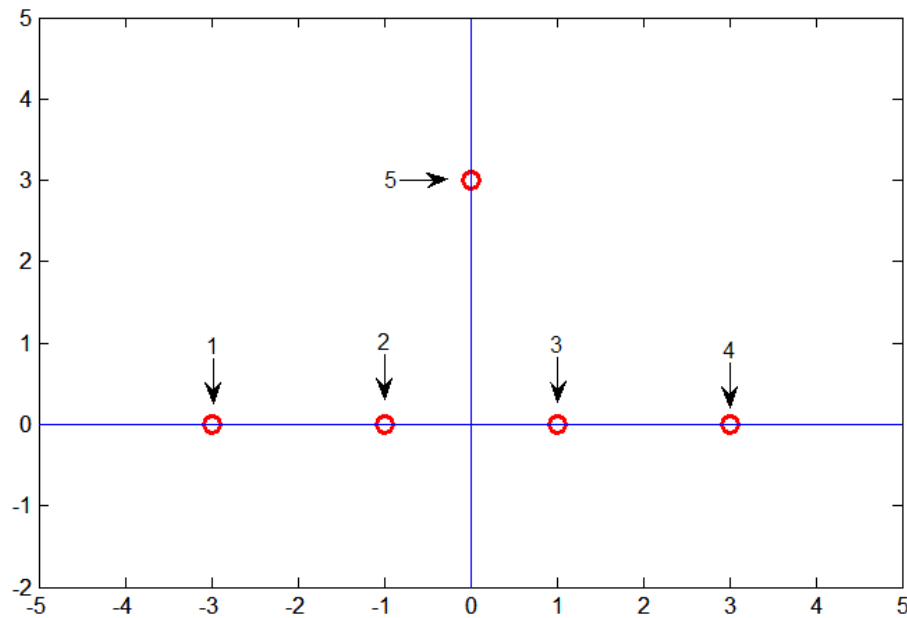
# Hierarchical Clustering: Centroid
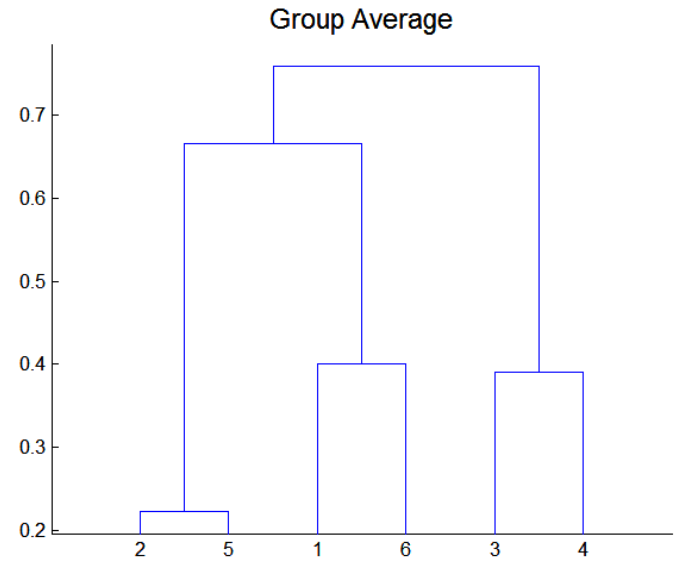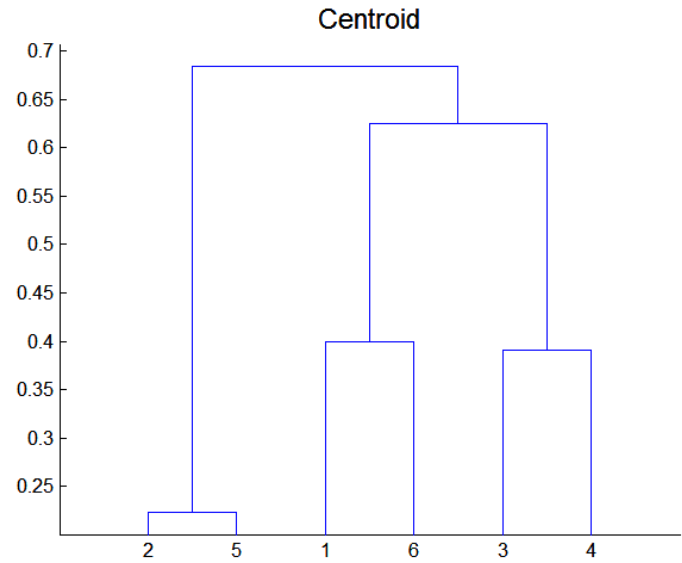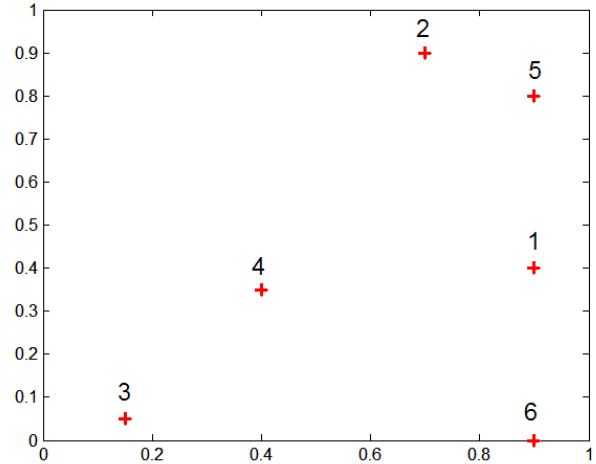


Nested Clusters

Dendrogram

# Inversion

- Two clusters that are merged may be more similar than the pair of clusters that were merged in a previous step
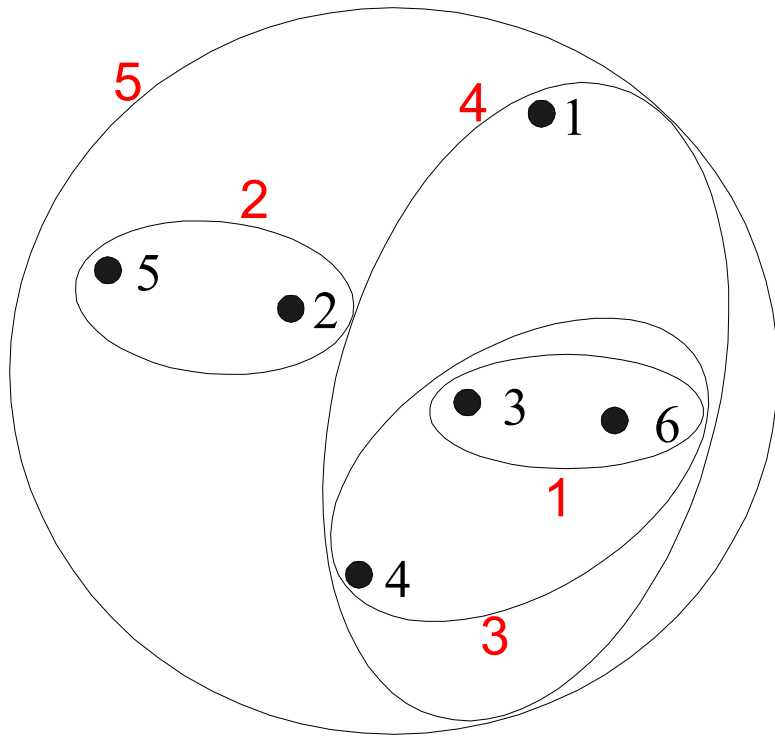
# Centroid vs Group Average

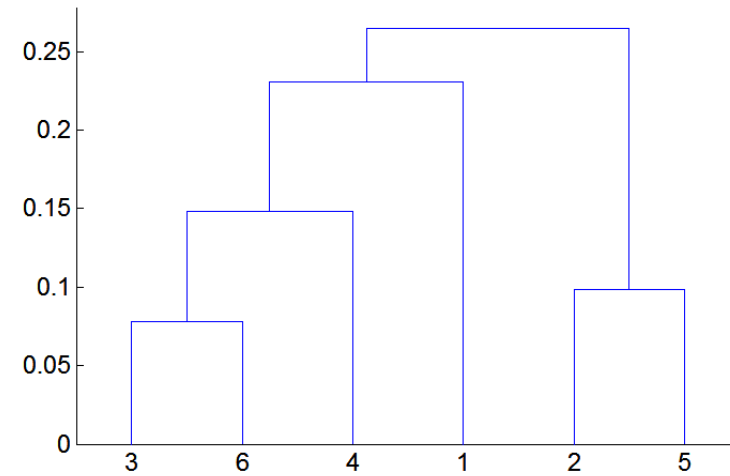| X | Y |
|---|---|
| 0.9 | 0.4 |
| 0.7 | 0.9 |
| 0.15 | 0.05 |
| 0.4 | 0.35 |
| 0.9 | 0.8 |
| 0.9 | 0 |

# Cluster Similarity: Ward's Method

- Proximity of two clusters is based on the increase in squared error when two clusters are merged
  - Similar to group average if distance between points is distance squared

- Less susceptible to noise

- Biased towards globular clusters
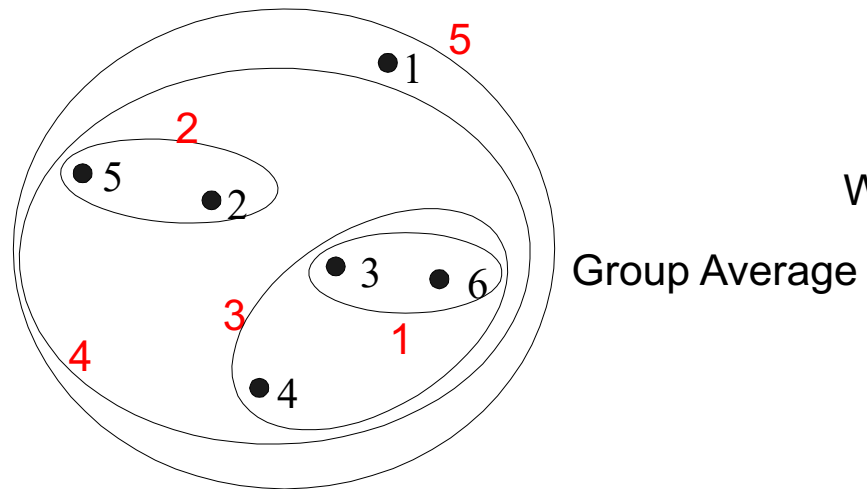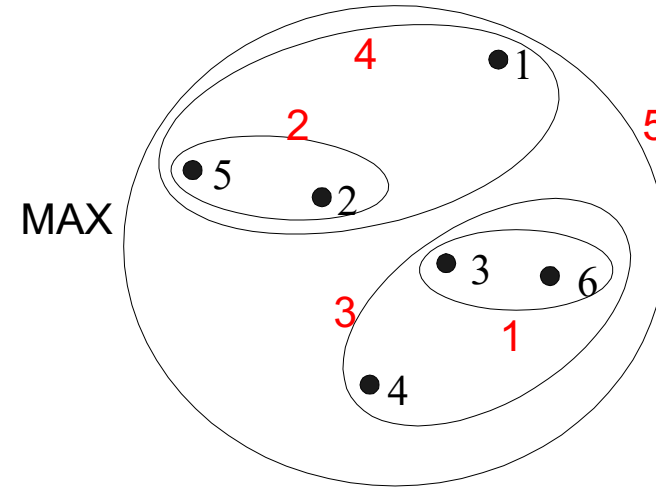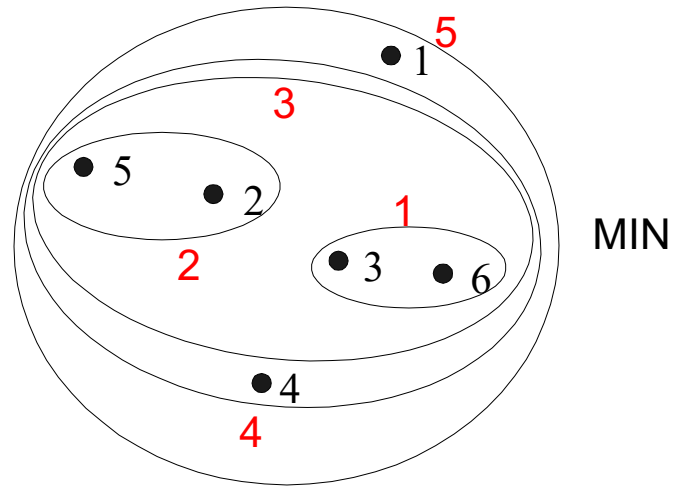
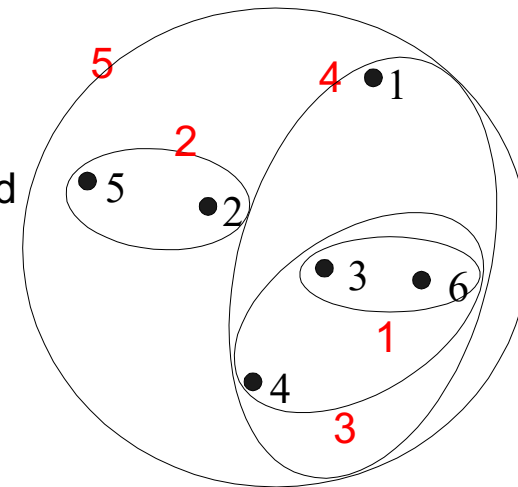# Hierarchical Clustering: Ward's Method



Nested Clusters

Dendrogram

# Hierarchical Clustering: Comparison

# Lance-Williams Formula

- Proximity between clusters Q and R, where R is formed by merging clusters A and B

$$p(R, Q) = \alpha_A p(A, Q) + \alpha_B p(B, Q) + \beta p(A, B) + \gamma |p(A, Q) - p(B, Q)|$$

| Clustering Method | $\alpha_{\mathbf{A}}$ | $\alpha_{\mathbf{B}}$ | $\beta$ | $\gamma$ |
|---|---|---|---|---|
| MIN | $1/2$ | $1/2$ | $0$ | $-1/2$ |
| MAX | $1/2$ | $1/2$ | $0$ | $1/2$ |
| Group Average | $\dfrac{m_A}{m_A + m_B}$ | $\dfrac{m_B}{m_A + m_B}$ | $0$ | $0$ |
| Centroid | $\dfrac{m_A}{m_A + m_B}$ | $\dfrac{m_B}{m_A + m_B}$ | $\dfrac{-m_A m_B}{(m_A + m_B)^2}$ | $0$ |
| Ward's | $\dfrac{m_A + m_Q}{m_A + m_B + m_Q}$ | $\dfrac{m_B + m_Q}{m_A + m_B + m_Q}$ | $\dfrac{-m_Q}{m_A + m_B + m_Q}$ | $0$ |

- $m_i$ is size of cluster i, p(i,j) is proximity of clusters i & j

# Hierarchical Clustering:  Problems and Limitations

- Once a decision is made to combine two clusters, it cannot be undone

- No global objective function is directly minimized

- Different schemes have problems with one or more of the following:
  - Sensitivity to noise
  - Difficulty handling different sized clusters and convex shapes
  - Breaking large clusters