# COM/BLM 453
# Data Mining
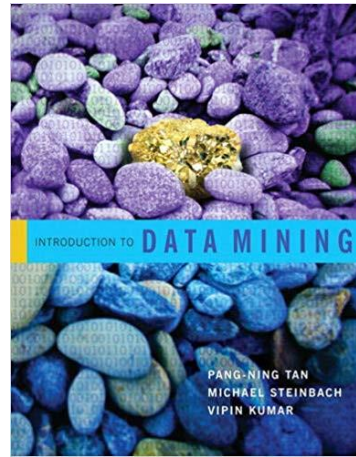## Asst. Prof. Dr. Bulent TUGRUL
### btugrul@eng.ankara.edu.tr

Slides are mainly based on:

Introduction to Data Mining

by Pang-Ning Tan, Michael Steinbach, Vipin Kumar

Pearson, 1st Edition, 2005

# Data Mining
# Cluster Analysis: Basic Concepts and Algorithms

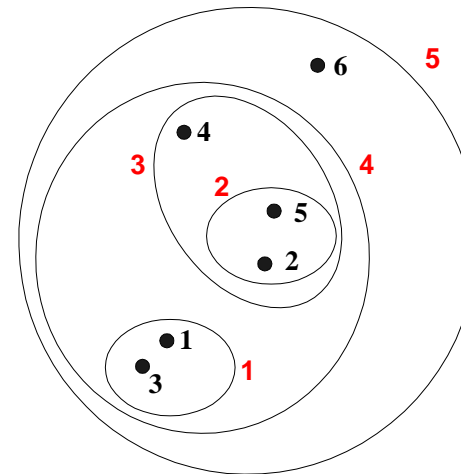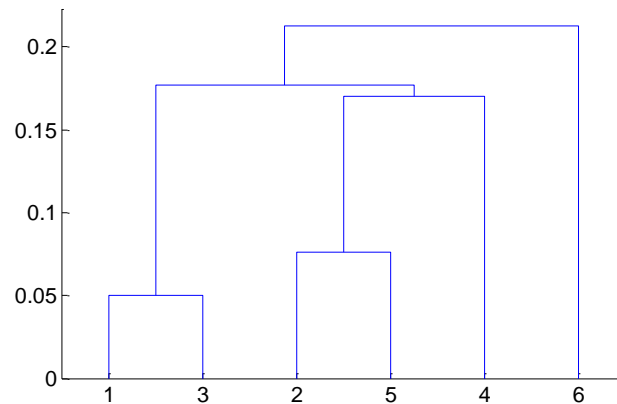## Lecture Notes for Chapter 8

## Introduction to Data Mining

by

Tan, Steinbach, Kumar

# 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 dendogram at the proper level

- They may correspond to meaningful taxonomies
  - Example in biological sciences (e.g., animal kingdom, phylogeny reconstruction, …)
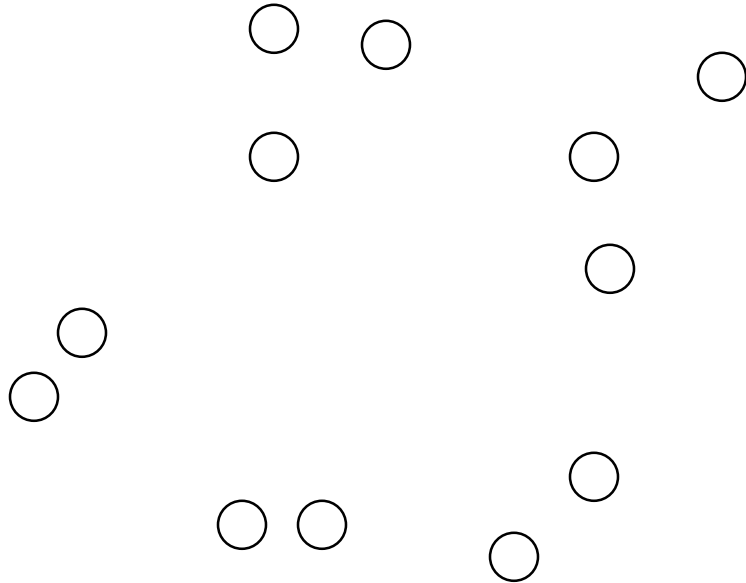
# Hierarchical Clustering

- 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

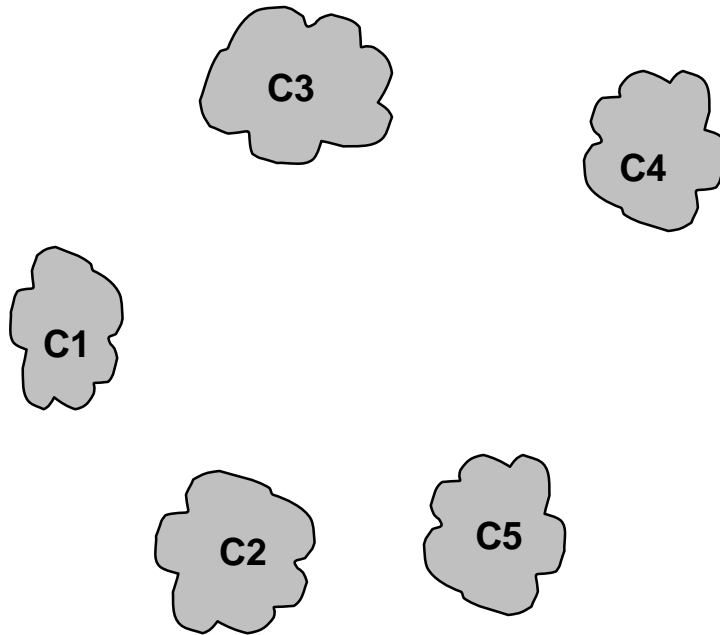- Start with clusters of individual points and a proximity matrix

|    | 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



**Proximity Matrix**

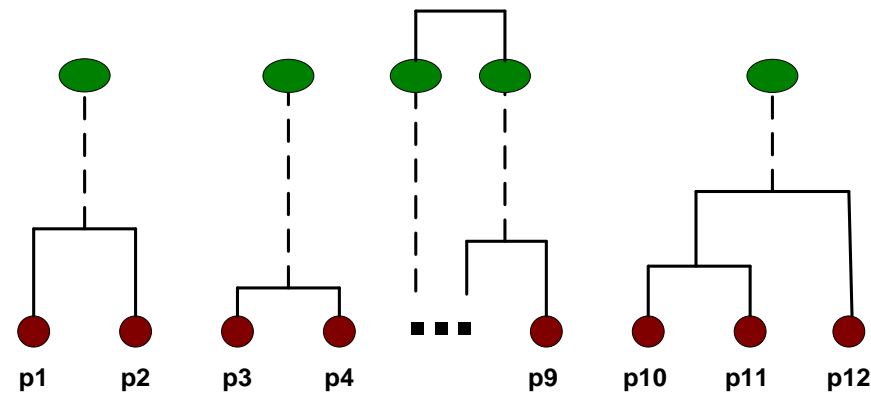|      | C1 | C2 | C3 | C4 | C5 |
|------|----|----|----|----|----|
| C1   |    |    |    |    |    |
| C2   |    |    |    |    |    |
| C3   |    |    |    |    |    |
| C4   |    |    |    |    |    |
| C5   |    |    |    |    |    |

# Intermediate Situation

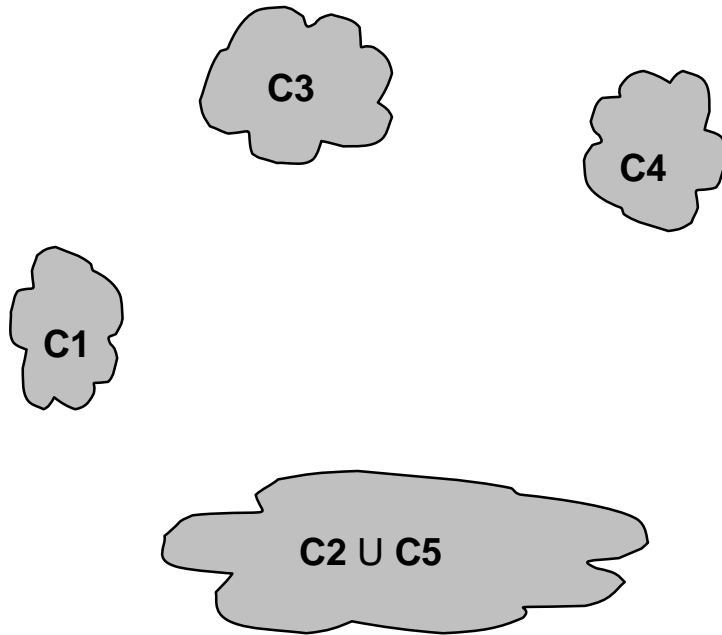- We want to merge the two closest clusters (C2 and C5) and update the proximity matrix.

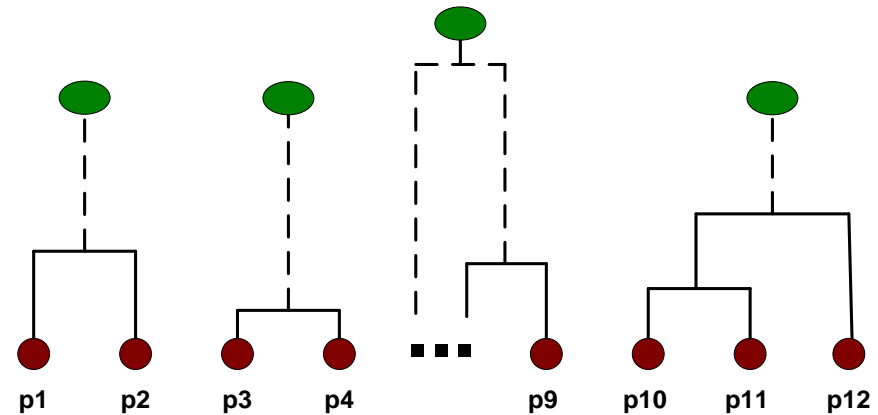|      | C1 | C2 | C3 | C4 | C5 |
|------|----|----|----|----|----|
| C1   |    |    |    |    |    |
| C2   |    |    |    |    |    |
| C3   |    |    |    |    |    |
| C4   |    |    |    |    |    |
| C5   |    |    |    |    |    |

**Proximity Matrix**

# After Merging

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

|          | C1 | C2 ∪ C5 | C3 | C4 |
|----------|----|---------|----|----|
| C1       |    | ?       |    |    |
| C2 ∪ C5  | ?  | ?       | ?  | ?  |
| C3       |    | ?       |    |    |
| C4       |    | ?       |    |    |

**Proximity Matrix**

# How to Define Inter-Cluster Similarity



Similarity?

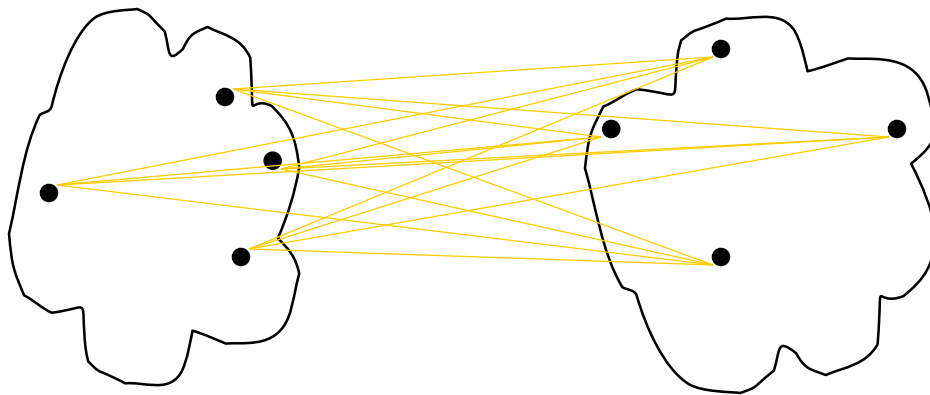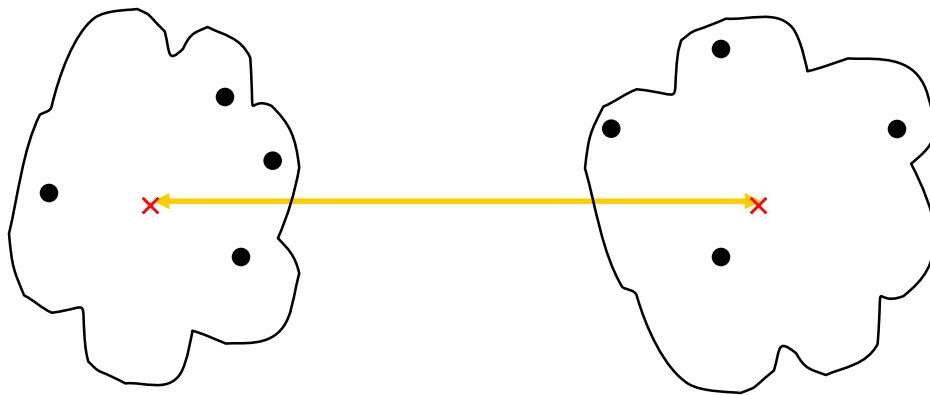|    | p1 | p2 | p3 | p4 | p5 | . . . |
|----|----|----|----|----|----|-------|
| p1 |    |    |    |    |    |       |
| p2 |    |    |    |    |    |       |
| p3 |    |    |    |    |    |       |
| p4 |    |    |    |    |    |       |
| p5 |    |    |    |    |    |       |
| .  |    |    |    |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

**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** | | | | | | |
| **.** | | | | | | |
| **.** | | | | | | |
| **.** | | | | | | |

**Proximity Matrix**

- <span style="color:red">MIN</span>
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
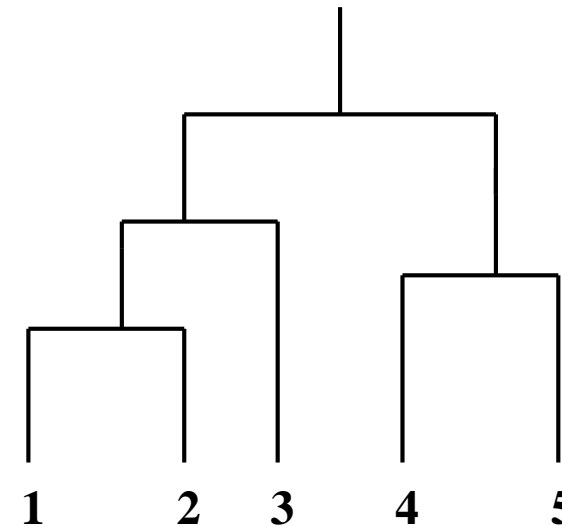  - Ward's Method uses squared error

# How to Define Inter-Cluster Similarity



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

**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** | | | | | | |
| . | | | | | | |
| . | | | | | | |
| . | | | | | | |

**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 | | | | | | |
| . | | | | | | |
| . | | | | | | |
| . | | | | | | |

**Proximity Matrix**

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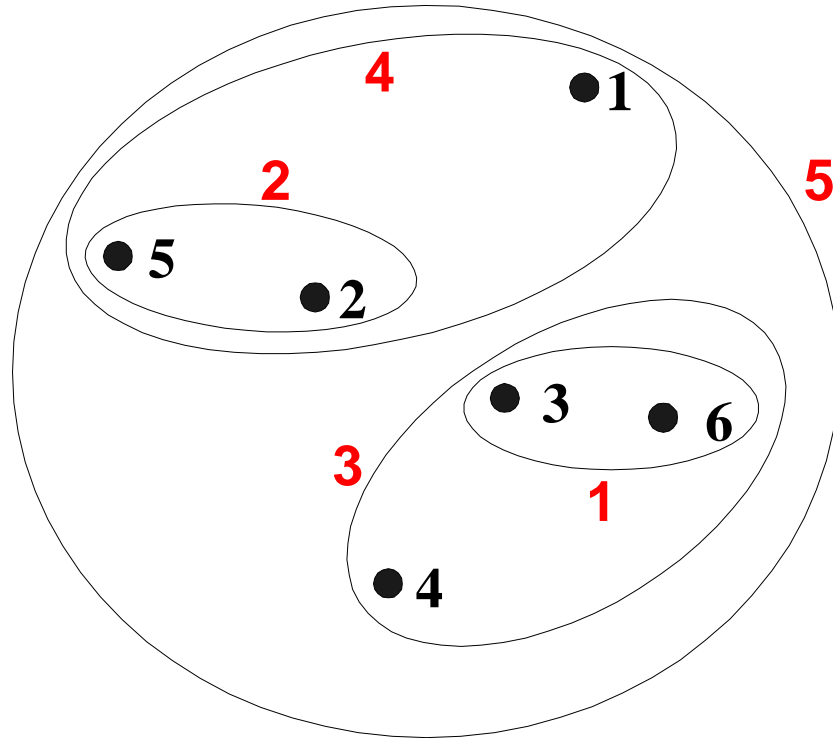- Other methods driven by an objective function
  - Ward's Method uses squared error

# 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 |

# Hierarchical Clustering: MIN



**Nested Clusters**

**Dendrogram**

# 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**

**Dendrogram**

# 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

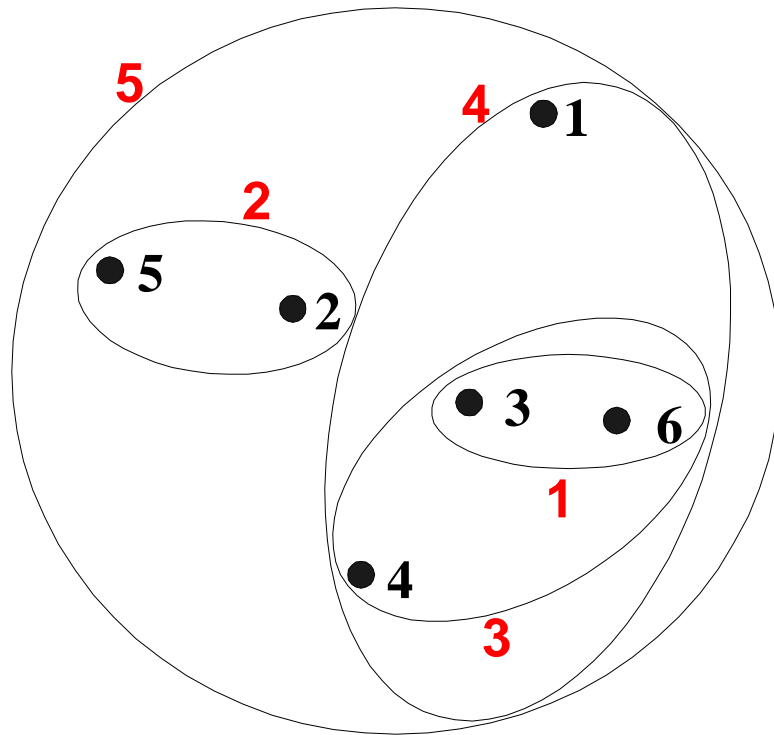- Proximity of two clusters is the average of pairwise proximity between points in the two clusters.

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

- 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
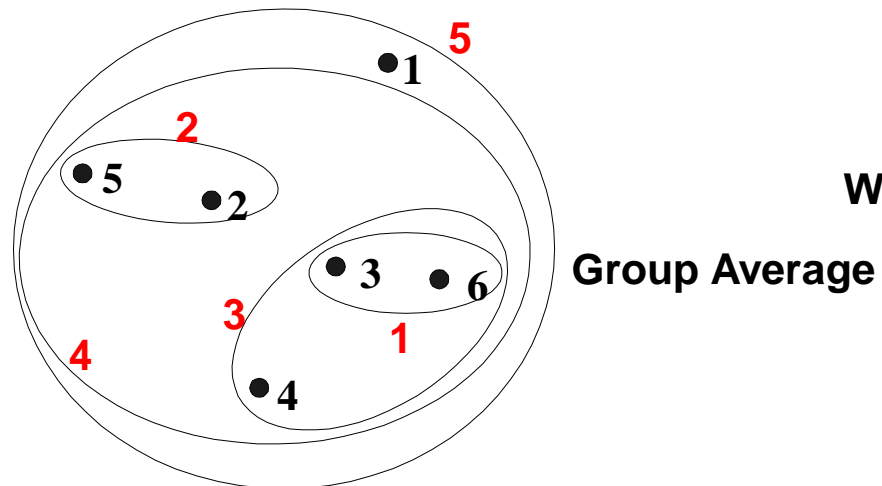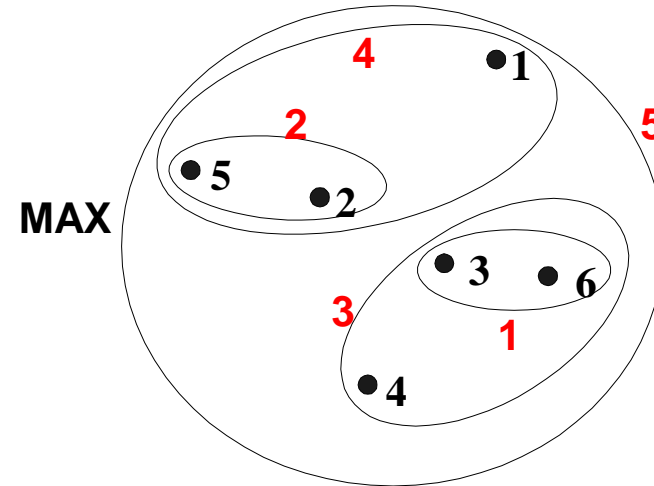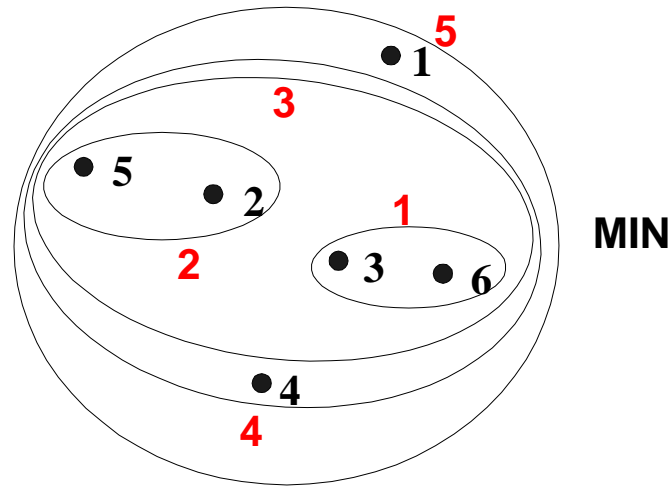


**Nested Clusters**

**Dendrogram**

# Hierarchical Clustering: Group Average

- Compromise between Single and Complete Link

- Strengths
  – Less susceptible to noise and outliers

- Limitations
  – Biased towards globular clusters

# Cluster Similarity: Ward's Method

- Similarity 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 and outliers

- Biased towards globular clusters

- Hierarchical analogue of K-means

  – Can be used to initialize K-means

# Hierarchical Clustering: Comparison



MIN

MAX

Group Average

Ward's Method

## 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

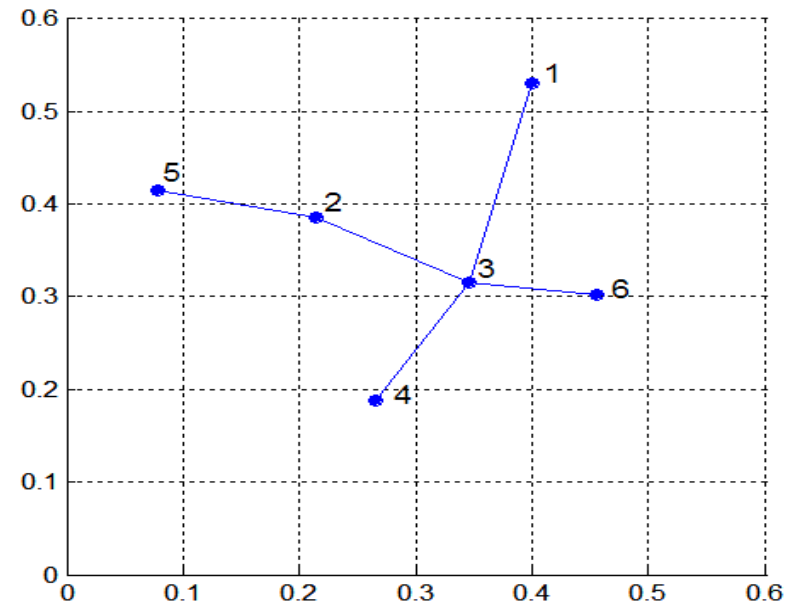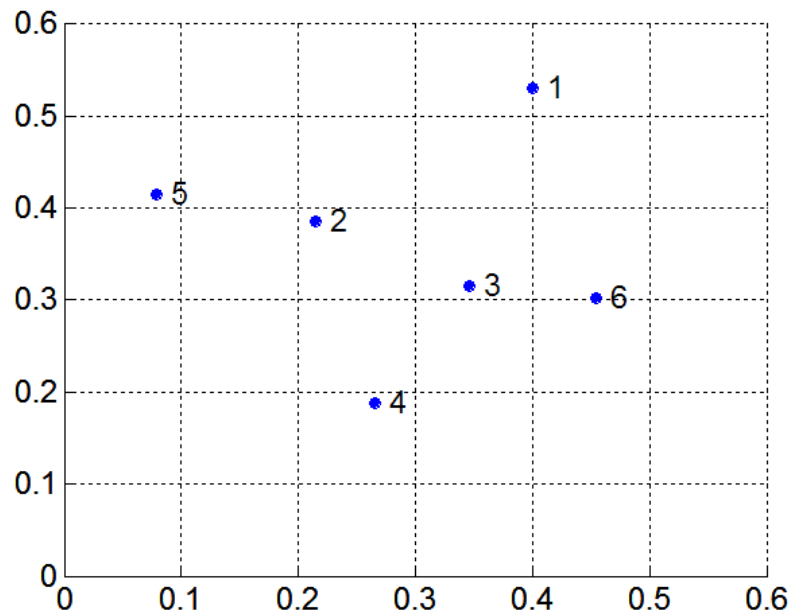  - Complexity can be reduced to O($N^2$ log(N) ) time for some approaches

# 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

# 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
  - Add q to the tree and put an edge between p and q

# 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

# DBSCAN

- DBSCAN is a density-based algorithm.

  – Density = number of points within a specified radius (Eps)

  – A point is a core point if it has more than a specified number of points (MinPts) within Eps
    - These are points that are at the interior of a cluster

  – A border point has fewer than MinPts within Eps, but is in the neighborhood of a core point

  – A noise point is any point that is not a core point or a border point.

# DBSCAN: Core, Border, and Noise Points

# DBSCAN Algorithm

- Eliminate noise points
- Perform clustering on the remaining points

$current\_cluster\_label \leftarrow 1$

**for** all core points **do**

    **if** the core point has no cluster label **then**

        $current\_cluster\_label \leftarrow current\_cluster\_label + 1$

        Label the current core point with cluster label $current\_cluster\_label$

    **end if**

    **for** all points in the $Eps$-neighborhood, except $i^{th}$ the point itself **do**

        **if** the point does not have a cluster label **then**

            Label the point with cluster label $current\_cluster\_label$

        **end if**

    **end for**

**end for**

# DBSCAN: Core, Border and Noise Points



**Original Points**

**Point types: <span style="color:green">core</span>, <span style="color:blue">border</span> and <span style="color:red">noise</span>**

**Eps = 10, MinPts = 4**

# When DBSCAN Works Well



**Original Points**

**Clusters**

- **Resistant to Noise**

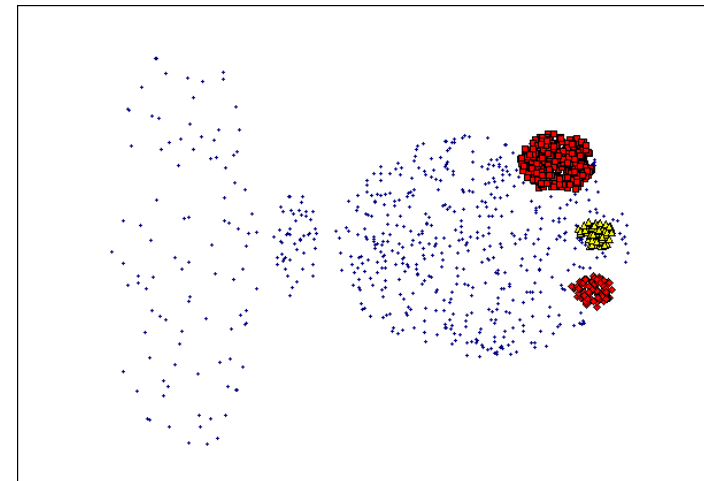- **Can handle clusters of different shapes and sizes**

# When DBSCAN Does NOT Work Well



**Original Points**

(MinPts=4, Eps=9.75).

(MinPts=4, Eps=9.92)

- **Varying densities**

- **High-dimensional data**

# Cluster Validity

- For supervised classification we have a variety of measures to evaluate how good our model is
  - Accuracy, precision, recall

- For cluster analysis, the analogous question is how to evaluate the "goodness" of the resulting clusters?

- But "clusters are in the eye of the beholder"!

- Then why do we want to evaluate them?
  - To avoid finding patterns in noise
  - To compare clustering algorithms
  - To compare two sets of clusters
  - To compare two clusters

# Clusters found in Random Data



Random Points

DBSCAN

K-means

Complete Link

# Measures of Cluster Validity

- Numerical measures that are applied to judge various aspects of cluster validity, are classified into the following three types.
    - External Index: Used to measure the extent to which cluster labels match externally supplied class labels.
        - Entropy
    - Internal Index: Used to measure the goodness of a clustering structure *without* respect to external information.
        - Sum of Squared Error (SSE)
    - Relative Index: Used to compare two different clusterings or clusters.
        - Often an external or internal index is used for this function, e.g., SSE or entropy
- Sometimes these are referred to as criteria instead of indices
    - However, sometimes criterion is the general strategy and index is the numerical measure that implements the criterion.
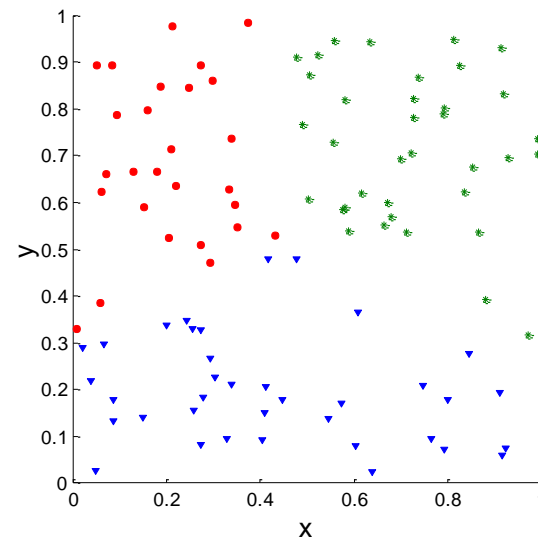
# Measuring Cluster Validity Via Correlation

- ## Two matrices

  - Proximity Matrix

  - "Incidence" Matrix

    - One row and one column for each data point

    - An entry is 1 if the associated pair of points belong to the same cluster

    - An entry is 0 if the associated pair of points belongs to different clusters

- ## Compute the correlation between the two matrices

  - Since the matrices are symmetric, only the correlation between $n(n-1) / 2$ entries needs to be calculated.

- ## High correlation indicates that points that belong to the same cluster are close to each other.

- ## Not a good measure for some density or contiguity based clusters.

# Measuring Cluster Validity Via Correlation

- Correlation of incidence and proximity matrices for the K-means clusterings of the following two data sets.



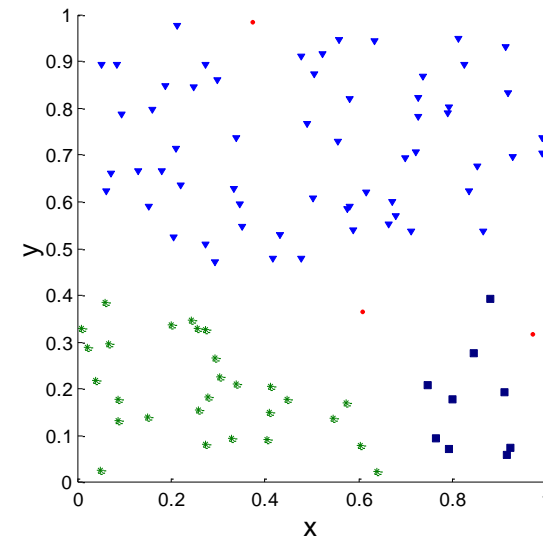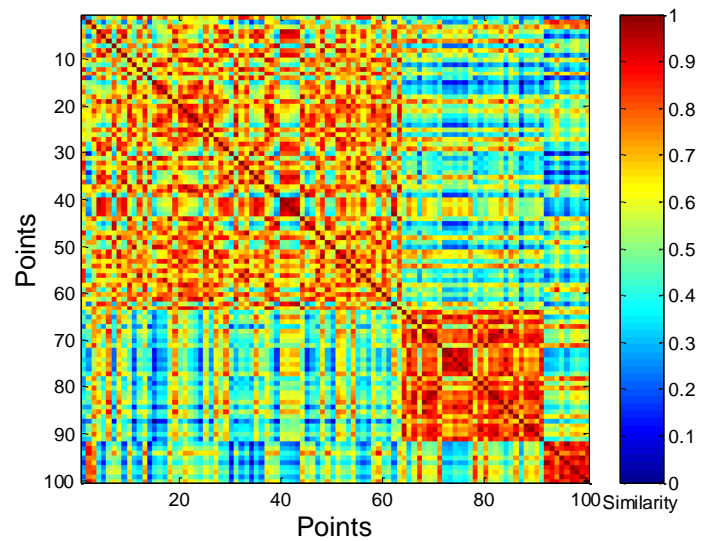Corr = -0.9235                    Corr = -0.5810

# Using Similarity Matrix for Cluster Validation

- Order the similarity matrix with respect to cluster labels and inspect visually.

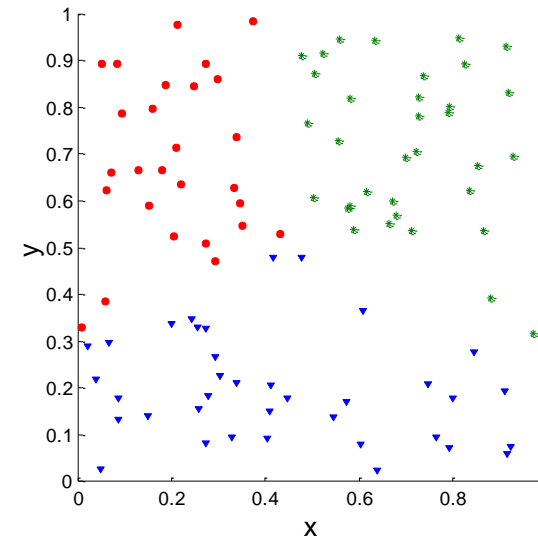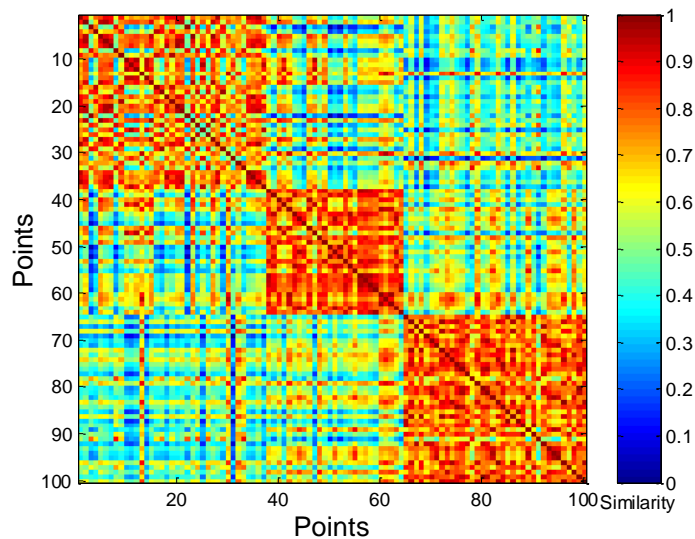# Using Similarity Matrix for Cluster Validation

- Clusters in random data are not so crisp



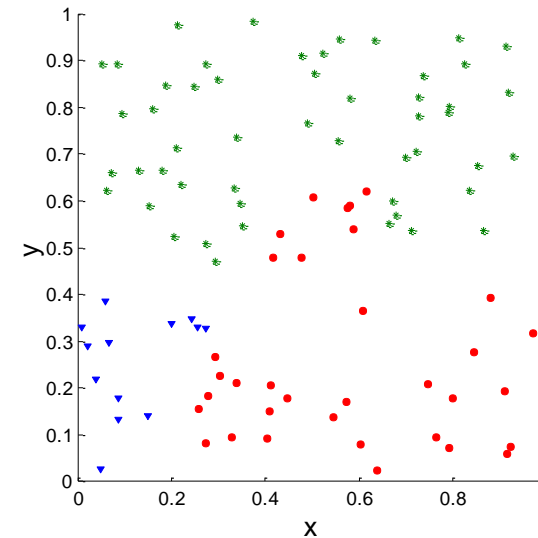**DBSCAN**
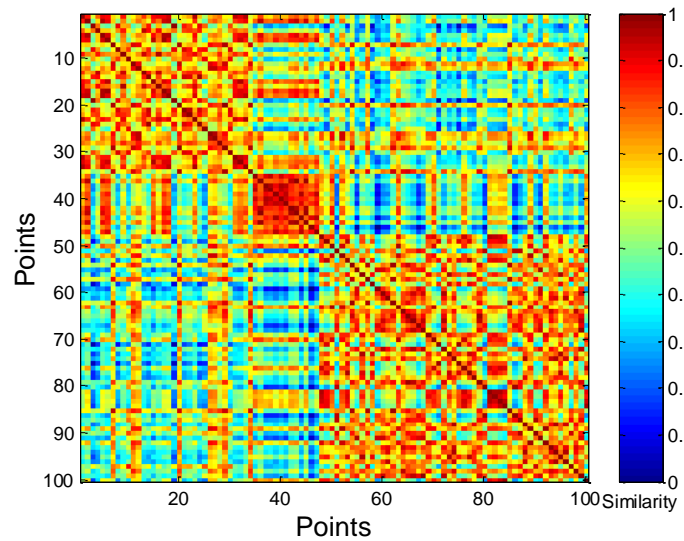
# Using Similarity Matrix for Cluster Validation

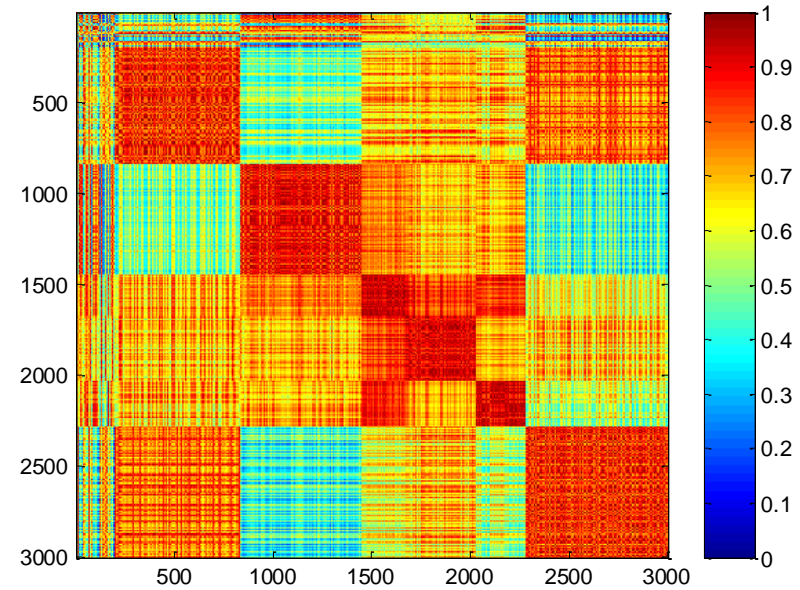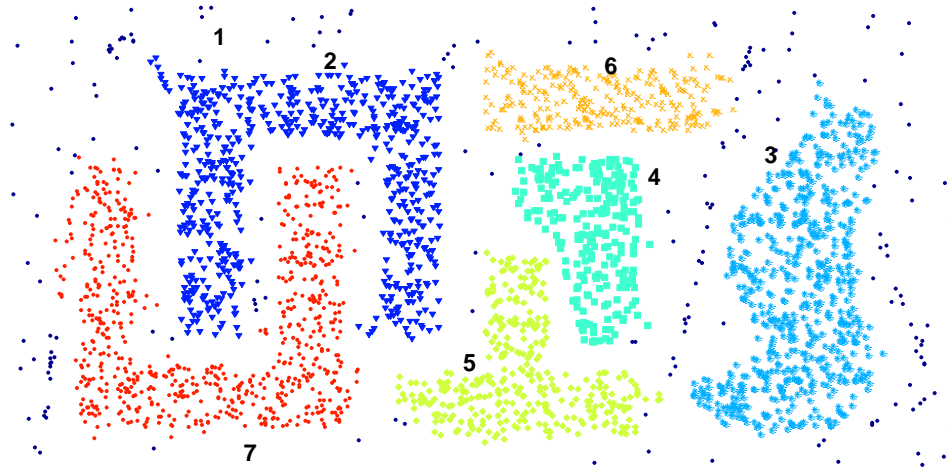- Clusters in random data are not so crisp



**K-means**

# Using Similarity Matrix for Cluster Validation

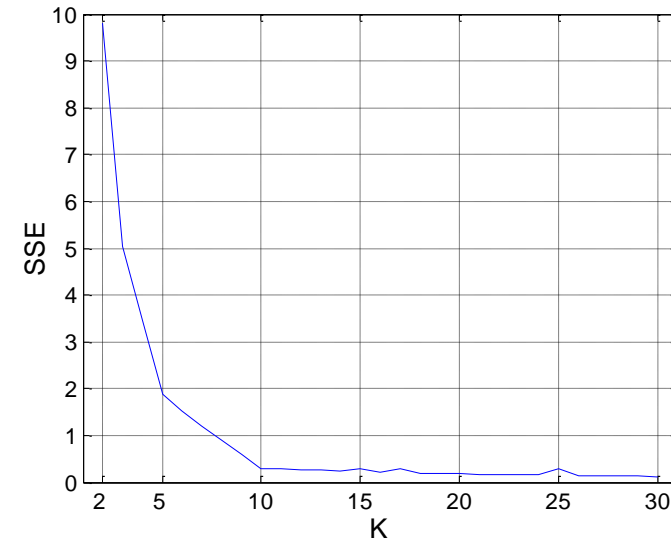- Clusters in random data are not so crisp
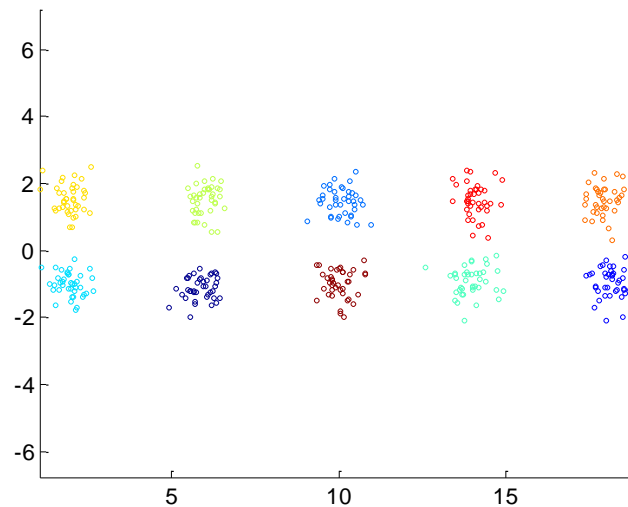


**Complete Link**

# Using Similarity Matrix for Cluster Validation
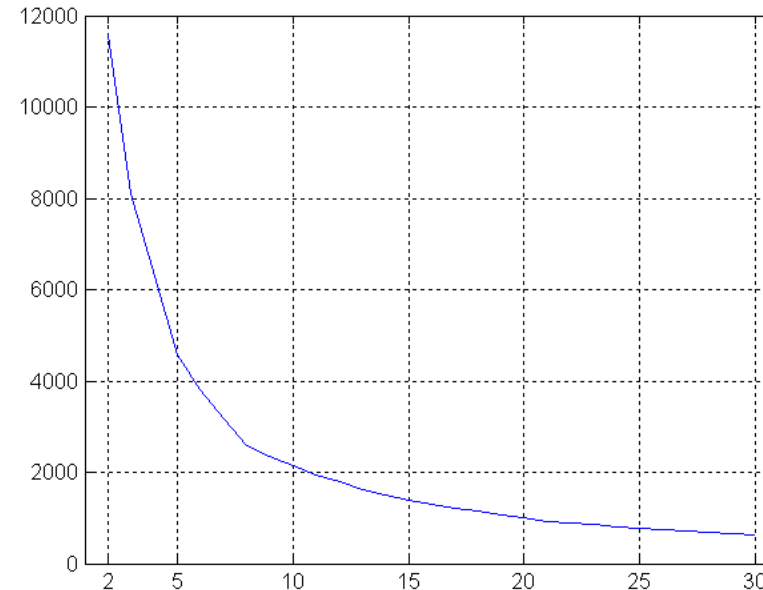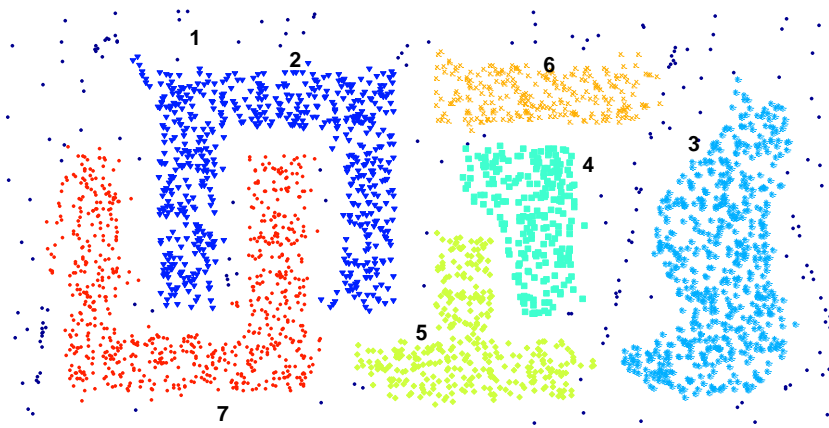


**DBSCAN**

# Internal Measures: SSE

- Clusters in more complicated figures aren't well separated
- Internal Index: Used to measure the goodness of a clustering structure without respect to external information
  - SSE
- SSE is good for comparing two clusterings or two clusters (average SSE).
- Can also be used to estimate the number of clusters

# Internal Measures: SSE

- SSE curve for a more complicated data set



**SSE of clusters found using K-means**

# Internal Measures: Cohesion and Separation

- **Cluster Cohesion:** Measures how closely related are objects in a cluster
  - Example: SSE
- **Cluster Separation:** Measure how distinct or well-separated a cluster is from other clusters
- Example: Squared Error
  - Cohesion is measured by the within cluster sum of squares (SSE)

$$WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$

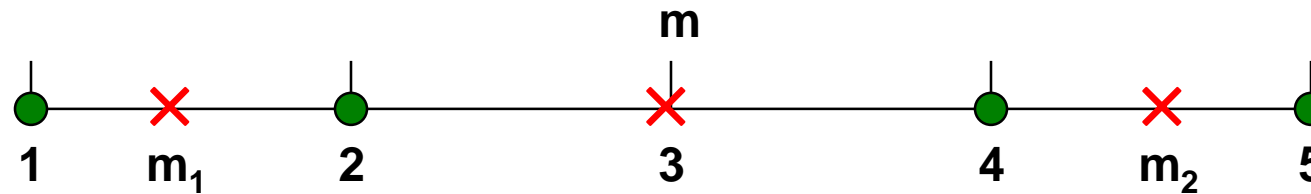  - Separation is measured by the between cluster sum of squares

$$BSS = \sum_i |C_i|(m - m_i)^2$$

  - Where $|C_i|$ is the size of cluster i

# Internal Measures: Cohesion and Separation

- Example: SSE
  - BSS + WSS = constant



**K=1 cluster:**

$$WSS = (1-3)^2 + (2-3)^2 + (4-3)^2 + (5-3)^2 = 10$$

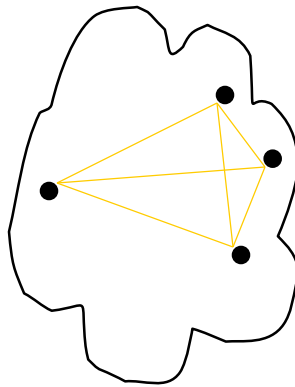$$BSS = 4 \times (3-3)^2 = 0$$

$$Total = 10 + 0 = 10$$

**K=2 clusters:**

$$WSS = (1-1.5)^2 + (2-1.5)^2 + (4-4.5)^2 + (5-4.5)^2 = 1$$

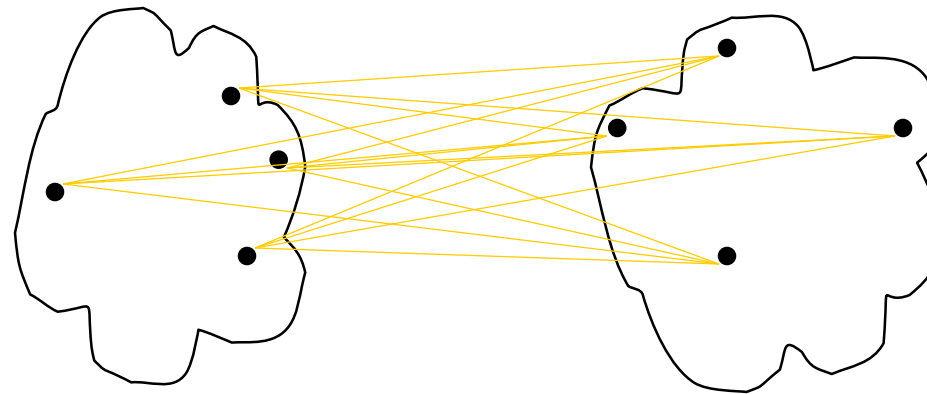$$BSS = 2 \times (3-1.5)^2 + 2 \times (4.5-3)^2 = 9$$

$$Total = 1 + 9 = 10$$

# Internal Measures: Cohesion and Separation

- A proximity graph based approach can also be used for cohesion and separation.

  – Cluster cohesion is the sum of the weight of all links within a cluster.

  – Cluster separation is the sum of the weights between nodes in the cluster and nodes outside the cluster.



cohesion                                    separation

# Final Comment on Cluster Validity

"The validation of clustering structures is the most difficult and frustrating part of cluster analysis.

Without a strong effort in this direction, cluster analysis will remain a black art accessible only to those true believers who have experience and great courage."

*Algorithms for Clustering Data*, Jain and Dubes