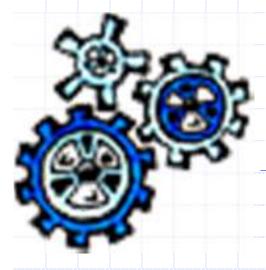# Data Mining - Clustering

Pisa KDD Lab, ISTI-CNR & Univ. Pisa
**http://www-kdd.isti.cnr.it/**

MAINS – Master in
Management dell'Innovazione
Scuola S. Anna
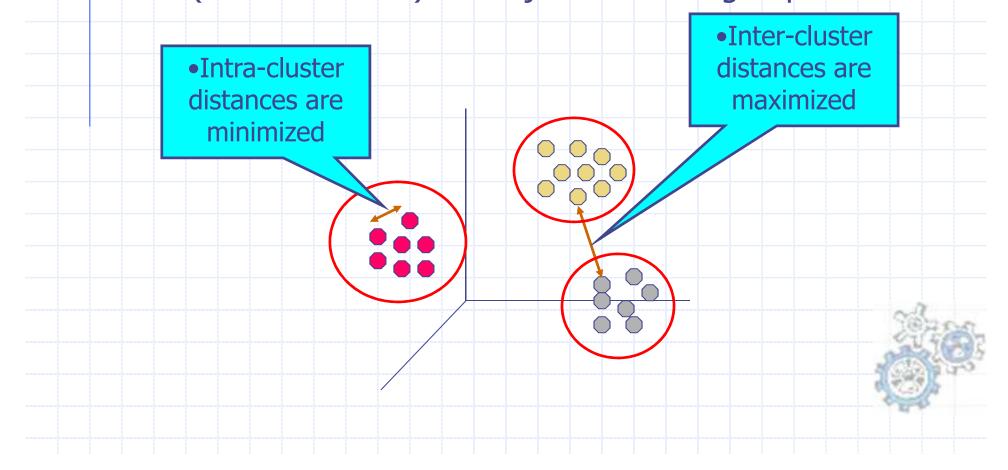
# Seminar 3 – Data Mining Technologies

# **Clustering**

# What is Cluster Analysis?

◆ Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups

•Intra-cluster distances are minimized

•Inter-cluster distances are maximized

# Applications of Cluster Analysis

◆ **Understanding**

- Group related documents for browsing, group genes and proteins that have similar functionality, or group stocks with similar price fluctuations

| | Discovered Clusters | Industry Group |
|---|---|---|
| 1 | Applied-Matl-DOWN,Bay-Network-Down,3-COM-DOWN, Cabletron-Sys-DOWN,CISCO-DOWN,HP-DOWN, DSC-Comm-DOWN,INTEL-DOWN,LSI-Logic-DOWN, Micron-Tech-DOWN,Texas-Inst-Down,Tellabs-Inc-Down, Natl-Semiconduct-DOWN,Oracl-DOWN,SGI-DOWN, Sun-DOWN | Technology1-DOWN |
| 2 | Apple-Comp-DOWN,Autodesk-DOWN,DEC-DOWN, ADV-Micro-Device-DOWN,Andrew-Corp-DOWN, Computer-Assoc-DOWN,Circuit-City-DOWN, Compaq-DOWN, EMC-Corp-DOWN, Gen-Inst-DOWN, Motorola-DOWN,Microsoft-DOWN,Scientific-Atl-DOWN | Technology2-DOWN |
| 3 | Fannie-Mae-DOWN,Fed-Home-Loan-DOWN, MBNA-Corp-DOWN,Morgan-Stanley-DOWN | Financial-DOWN |
| 4 | Baker-Hughes-UP,Dresser-Inds-UP,Halliburton-HLD-UP, Louisiana-Land-UP,Phillips-Petro-UP,Unocal-UP, Schlumberger-UP | Oil-UP |

◆ **Summarization**

- Reduce the size of large data sets



- Clustering precipitation in Australia

# What is not Cluster Analysis?

◆ **Supervised classification**
- Have class label information
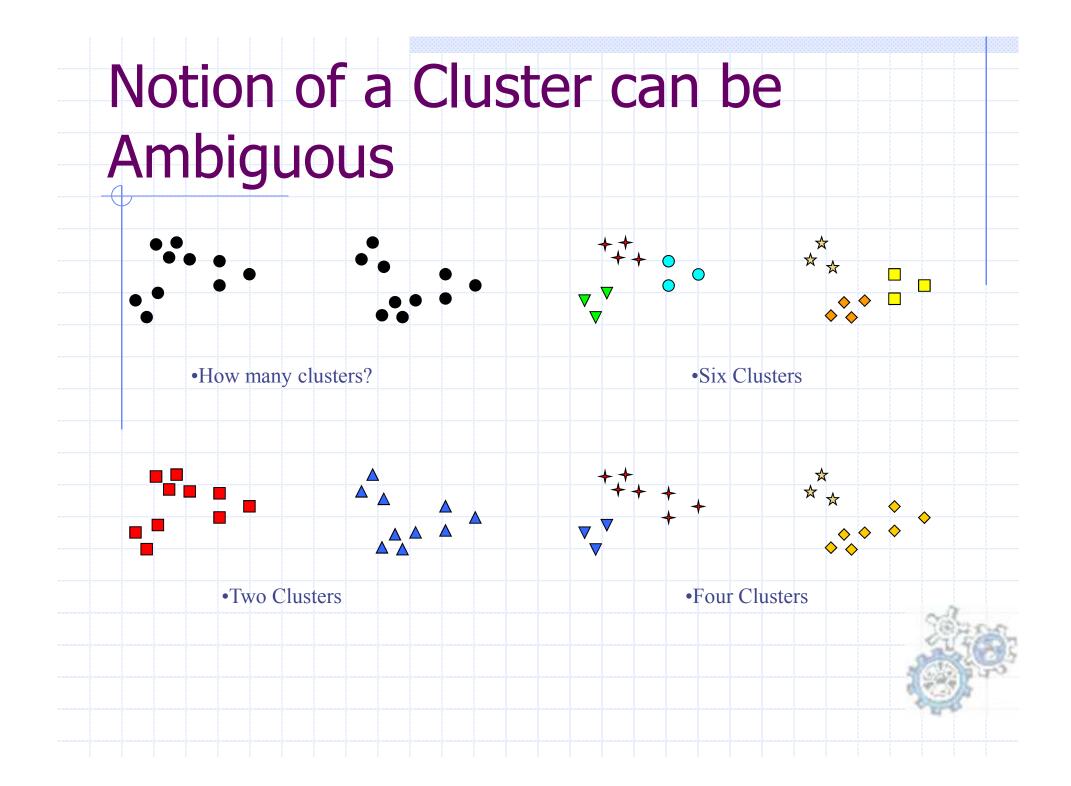
◆ **Simple segmentation**
- Dividing students into different registration groups alphabetically, by last name

◆ **Results of a query**
- Groupings are a result of an external specification

◆ **Graph partitioning**
- Some mutual relevance and synergy, but areas are not identical

# Notion of a Cluster can be Ambiguous

•How many clusters?

•Six Clusters

•Two Clusters

•Four Clusters

# Similarity and Dissimilarity

- ◈ Similarity
  - ▪ Numerical measure of how alike two data objects are.
  - ▪ Is higher when objects are more alike.
  - ▪ Often falls in the range [0,1]
- ◈ Dissimilarity
  - ▪ Numerical measure of how different are two data objects
  - ▪ Lower when objects are more alike
  - ▪ Minimum dissimilarity is often 0
  - ▪ Upper limit varies
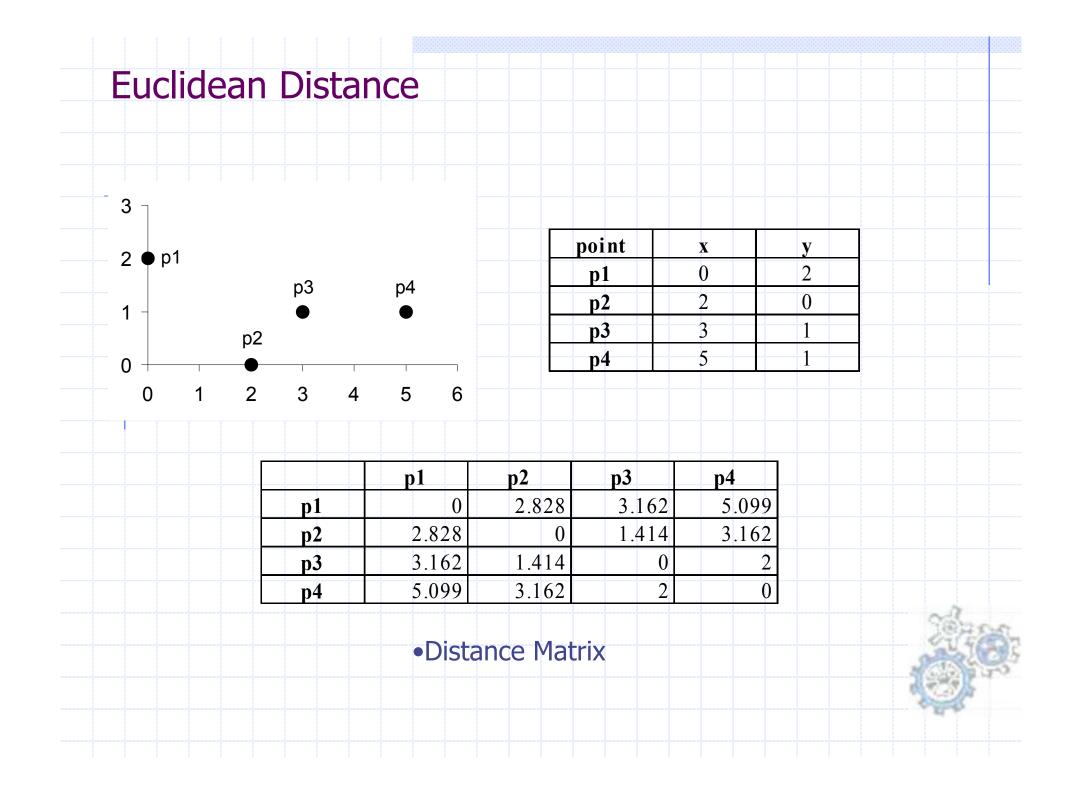- ◈ Proximity refers to a similarity or dissimilarity

# Euclidean Distance

◆ Euclidean Distance

$$dist = \sqrt{\sum_{k=1}^{n}(p_k - q_k)^2}$$

Where $n$ is the number of dimensions (attributes) and $p_k$ and $q_k$ are, respectively, the k[th] attributes (components) or data objects $p$ and $q$.

◆ Standardization is necessary, if scales differ.

# Euclidean Distance

| point | x | y |
|-------|---|---|
| p1 | 0 | 2 |
| p2 | 2 | 0 |
| p3 | 3 | 1 |
| p4 | 5 | 1 |

|  | p1 | p2 | p3 | p4 |
|------|------|------|------|------|
| p1 | 0 | 2.828 | 3.162 | 5.099 |
| p2 | 2.828 | 0 | 1.414 | 3.162 |
| p3 | 3.162 | 1.414 | 0 | 2 |
| p4 | 5.099 | 3.162 | 2 | 0 |

- Distance Matrix

# Similarity Between Binary Vectors

◆ Common situation is that objects, $p$ and $q$, have only binary attributes

◆ Compute similarities using the following quantities

$M_{01}$ = the number of attributes where p was 0 and q was 1
$M_{10}$ = the number of attributes where p was 1 and q was 0
$M_{00}$ = the number of attributes where p was 0 and q was 0
$M_{11}$ = the number of attributes where p was 1 and q was 1

◆ Jaccard Coefficient

J = number of 11 matches / number of not-both-zero attributes values

$$= (M_{11}) / (M_{01} + M_{10} + M_{11})$$

# Jaccard: Example

$p =$  1 0 0 0 0 0 0 0 0 0

$q =$  0 0 0 0 0 0 1 0 0 1

$M_{01} = 2$  (the number of attributes where p was 0 and q was 1)
$M_{10} = 1$  (the number of attributes where p was 1 and q was 0)
$M_{00} = 7$  (the number of attributes where p was 0 and q was 0)
$M_{11} = 0$  (the number of attributes where p was 1 and q was 1)

$$J = (M_{11}) / (M_{01} + M_{10} + M_{11}) = 0 / (2 + 1 + 0) = 0$$

# Types of Clusterings

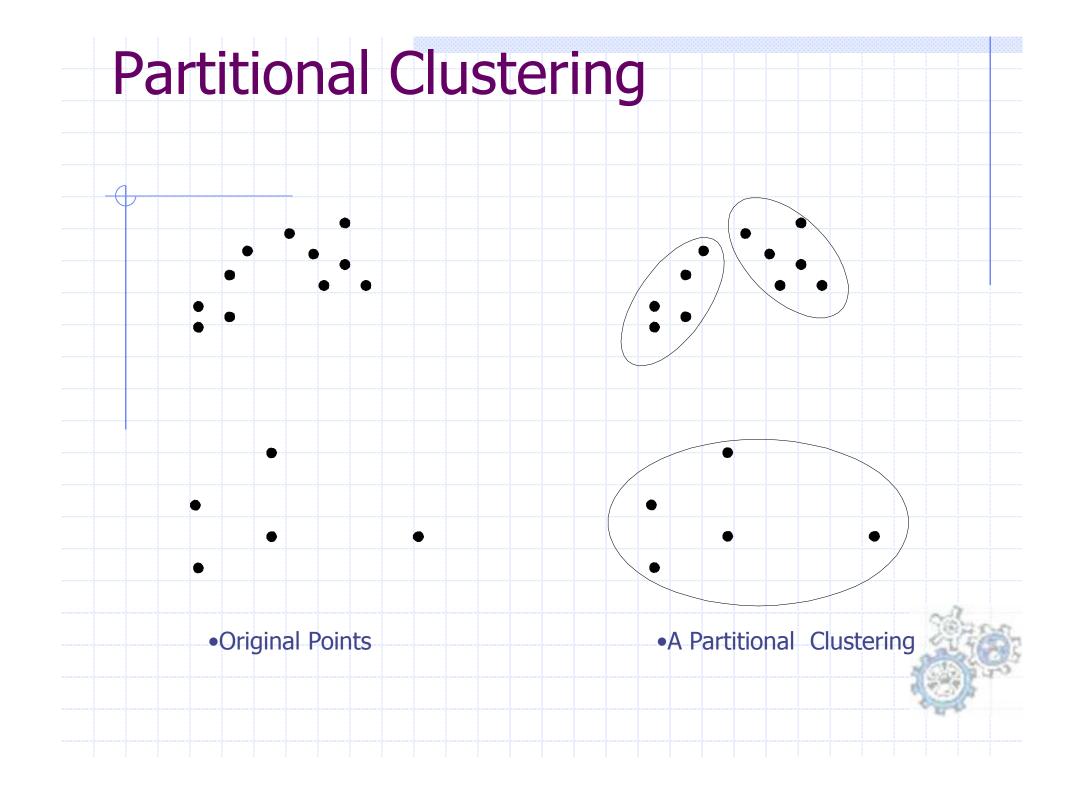◆ A clustering is a set of clusters

◆ Important distinction between hierarchical and partitional sets of clusters
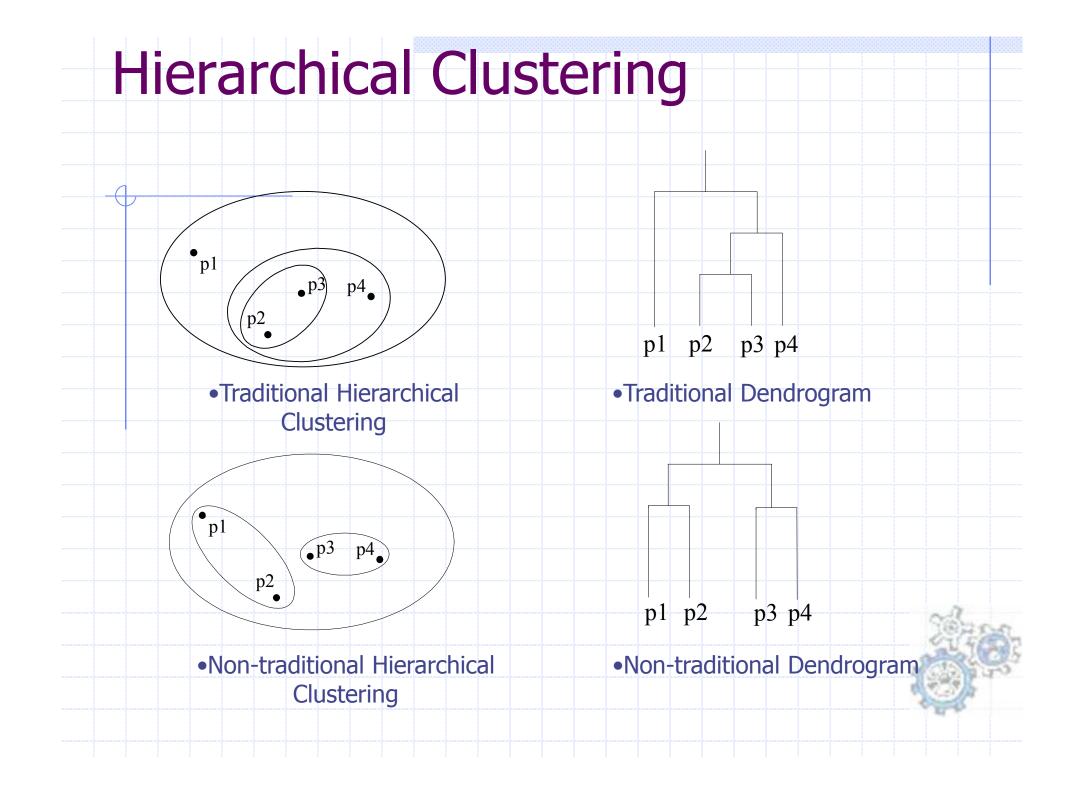
◆ Partitional Clustering
  - A division data objects into non-overlapping subsets (clusters) such that each data object is in exactly one subset

◆ Hierarchical clustering
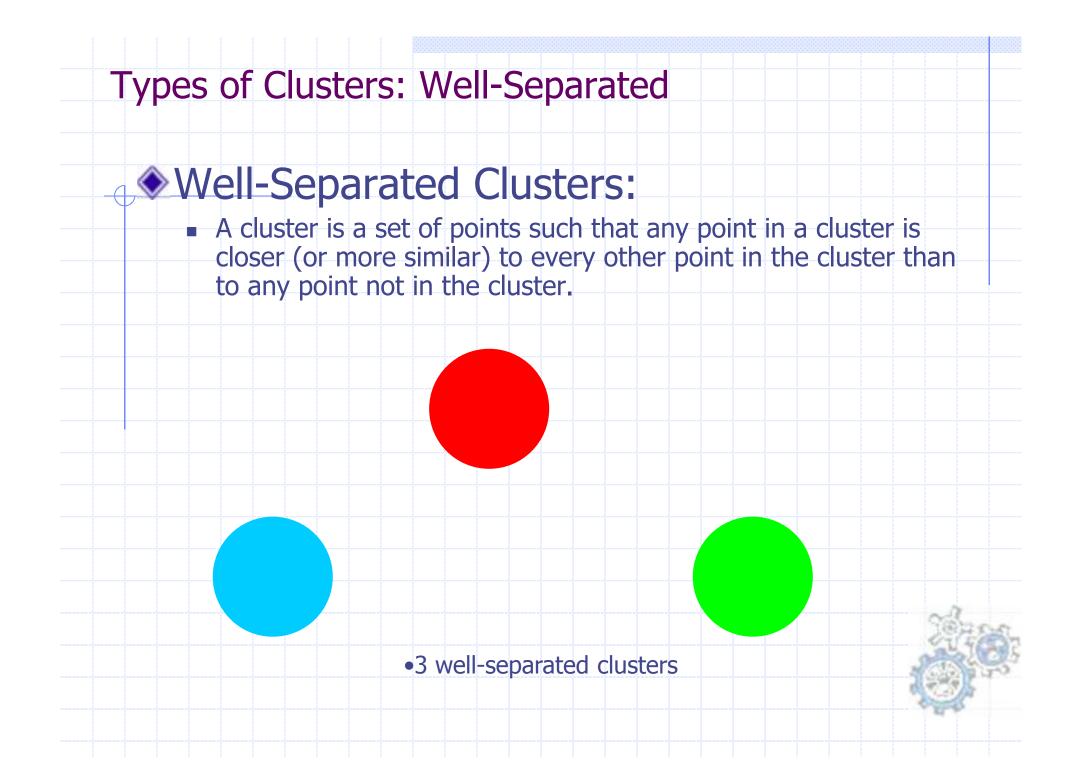  - A set of nested clusters organized as a hierarchical tree

# Partitional Clustering



• Original Points

• A Partitional Clustering

# Hierarchical Clustering

- Traditional Hierarchical Clustering

- Traditional Dendrogram

- Non-traditional Hierarchical Clustering

- Non-traditional Dendrogram

# Types of Clusters

- ◆ Well-separated clusters

- ◆ Center-based clusters

- ◆ Contiguous clusters

- ◆ Density-based clusters

- ◆ Property or Conceptual

## Well-Separated Clusters:

- A cluster is a set of points such that any point in a cluster is closer (or more similar) to every other point in the cluster than to any point not in the cluster.
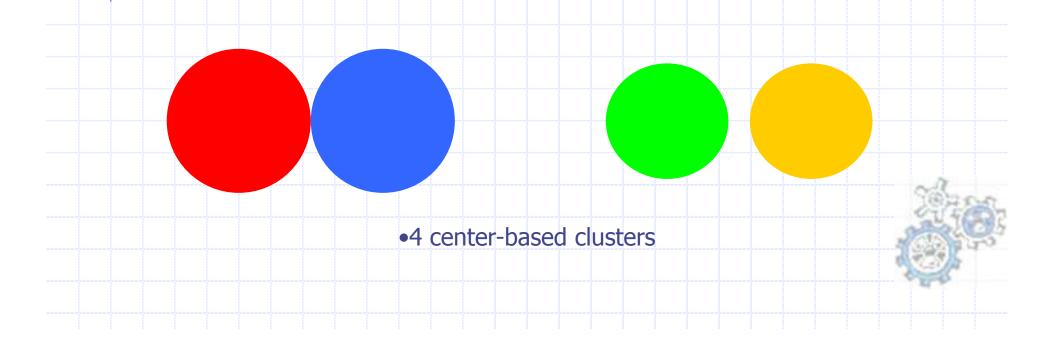
•3 well-separated clusters

# Types of Clusters: Center-Based
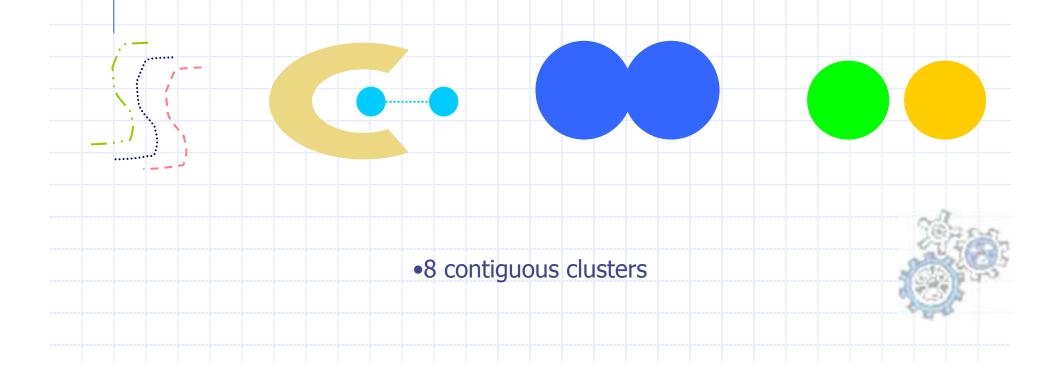
## ◆ Center-based

- A cluster is a set of objects such that an object in a cluster is closer (more similar) to the "center" of a cluster, than to the center of any other cluster

- The center of a cluster is often a centroid, the average of all the points in the cluster, or a medoid, the most "representative" point of a cluster

• 4 center-based clusters

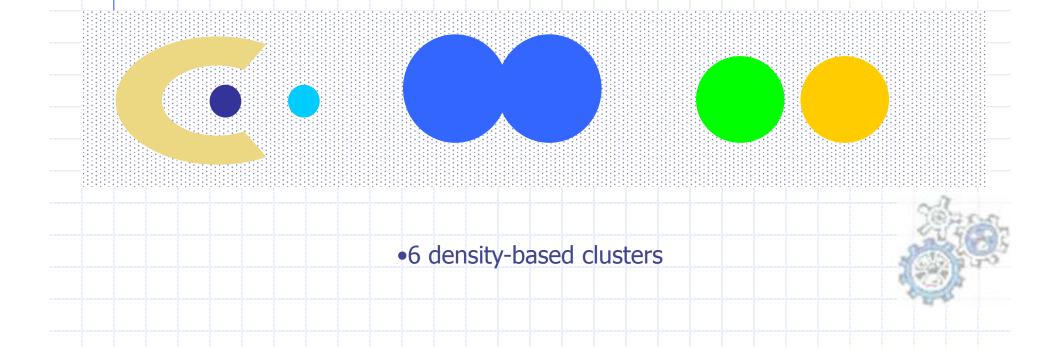◆ **Contiguous Cluster (Nearest neighbor or Transitive)**

- A cluster is a set of points such that a point in a cluster is closer (or more similar) to one or more other points in the cluster than to any point not in the cluster.

•8 contiguous clusters
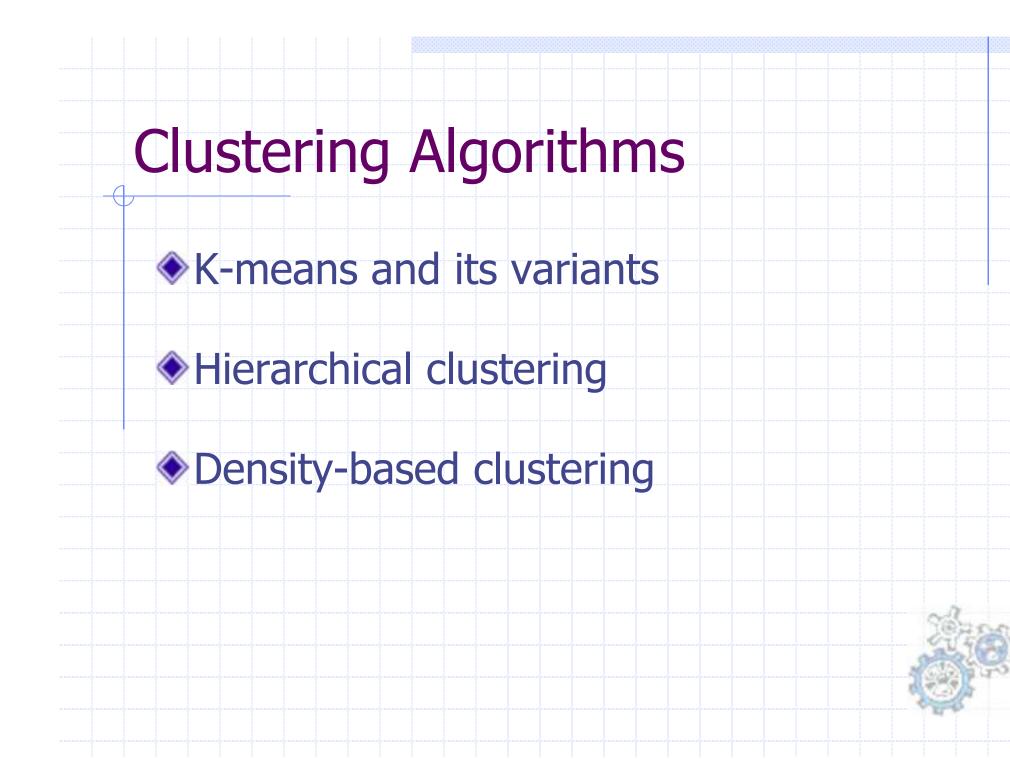
# Types of Clusters: Density-Based

## ◆ Density-based

- A cluster is a dense region of points, which is separated by low-density regions, from other regions of high density.
- Used when the clusters are irregular or intertwined, and when noise and outliers are present.



• 6 density-based clusters

# Characteristics of the Input Data Are Important

- ◆ **Type of proximity or density measure**
  - ■ This is a derived measure, but central to clustering
- ◆ **Sparseness**
  - ■ Dictates type of similarity
  - ■ Adds to efficiency
- ◆ **Attribute type**
  - ■ Dictates type of similarity
- ◆ **Type of Data**
  - ■ Dictates type of similarity
  - ■ Other characteristics, e.g., autocorrelation
- ◆ **Dimensionality**
- ◆ **Noise and Outliers**
- ◆ **Type of Distribution**

# Clustering Algorithms

◆ K-means and its variants

◆ Hierarchical clustering
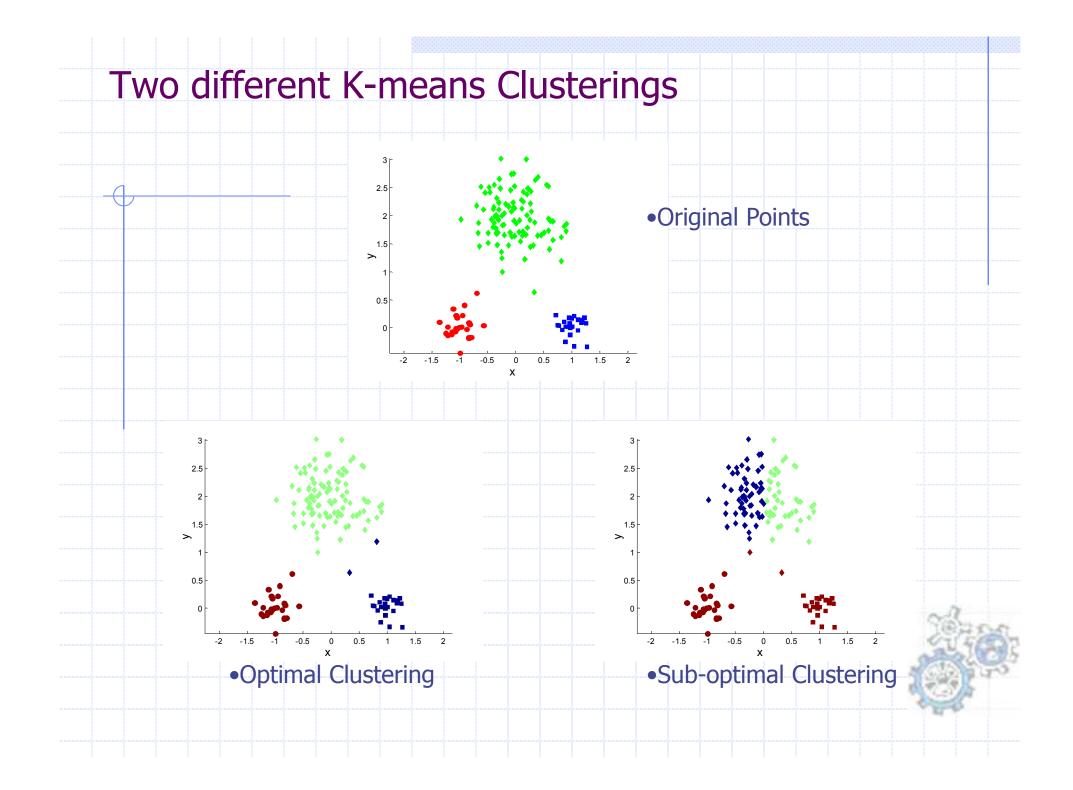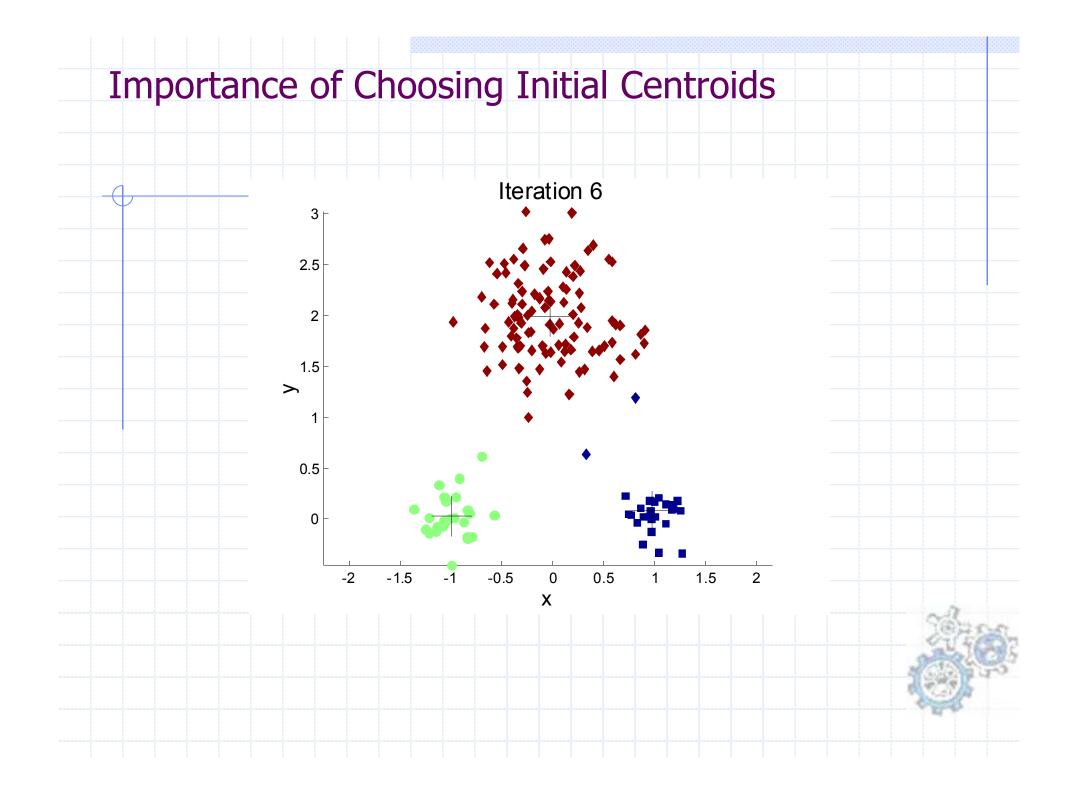
◆ Density-based clustering

# K-means Clustering

◆ Partitional clustering approach

◆ Each cluster is associated with a centroid (center point)

◆ Each point is assigned to the cluster with the closest centroid

◆ Number of clusters, K, must be specified

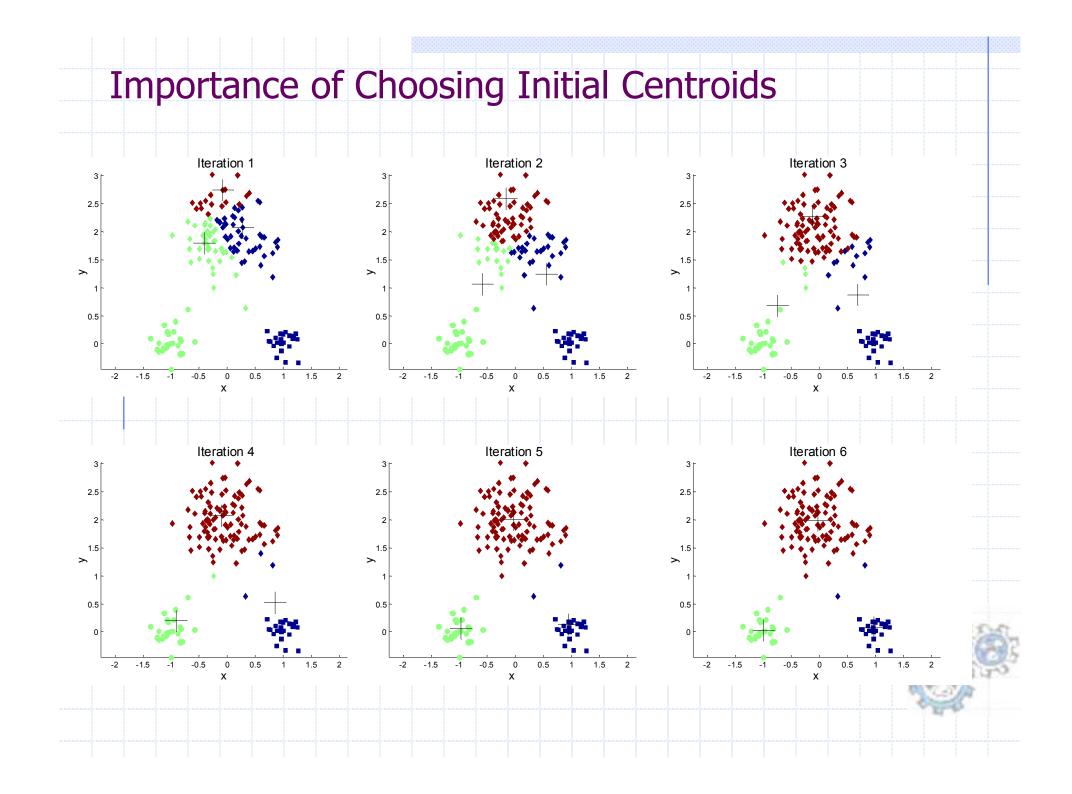◆ The basic algorithm is very simple

1: Select $K$ points as the initial centroids.

2: **repeat**

3:     Form $K$ clusters by assigning all points to the closest centroid.

4:     Recompute the centroid of each cluster.

5: **until** The centroids don't change
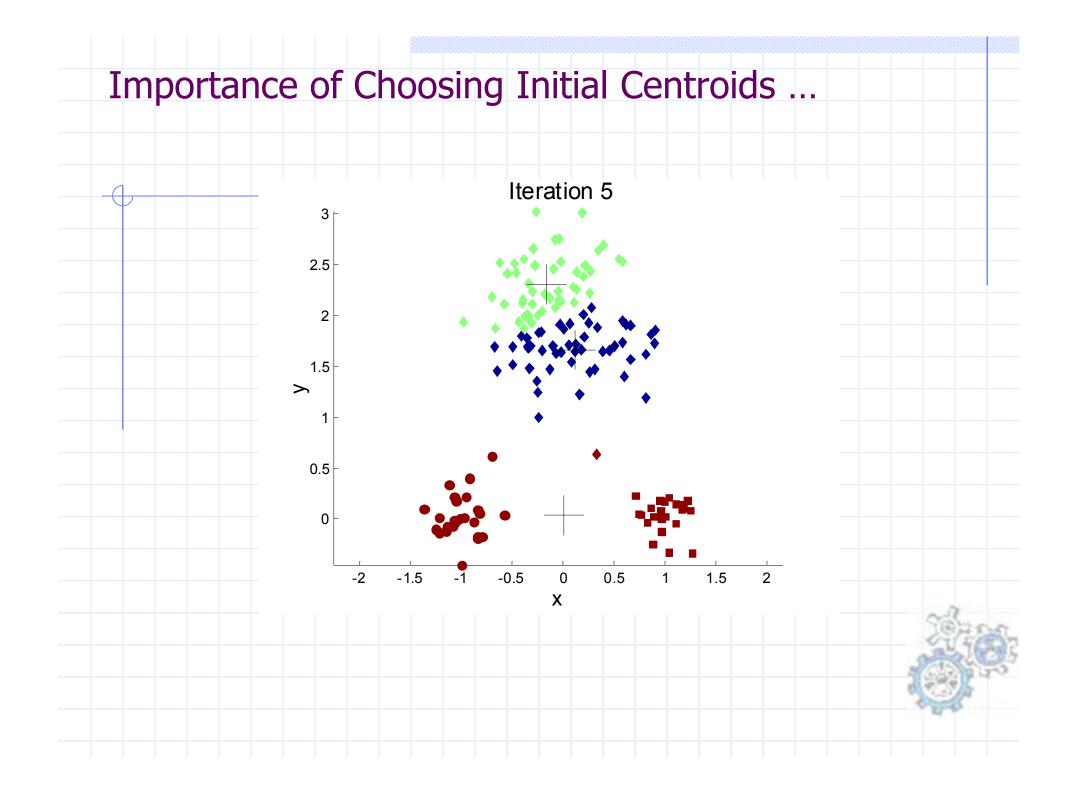
# K-means Clustering – Details

- ◆ Initial centroids are often chosen randomly.
  - ▪ ~~Clusters~~ produced vary from one run to another.
- ◆ The centroid is (typically) the mean of the points in the cluster.
- ◆ 'Closeness' is measured by Euclidean distance, cosine similarity, correlation, etc.
- ◆ K-means will converge for common similarity measures mentioned above.
- ◆ Most of the convergence happens in the first few iterations.
  - ▪ Often the stopping condition is changed to 'Until relatively few points change clusters'
- ◆ Complexity is O( n * K * I * d )
  - ▪ n = number of points, K = number of clusters, I = number of iterations, d = number of attributes

# Two different K-means Clusterings



- Original Points

- Optimal Clustering

- Sub-optimal Clustering

# Importance of Choosing Initial Centroids



Iteration 6

# Importance of Choosing Initial Centroids

# Importance of Choosing Initial Centroids …



Iteration 5

# Importance of Choosing Initial Centroids ...

# 10 Clusters Example

Iteration 4



• Starting with two initial centroids in one cluster of each pair of clusters

# 10 Clusters Example



• Starting with two initial centroids in one cluster of each pair of clusters

# 10 Clusters Example

Iteration 4



•Starting with some pairs of clusters having three initial centroids, while other have only one.

# 10 Clusters Example



• Starting with some pairs of clusters having three initial centroids, while other have only one.

# Limitations of K-means

◆ K-means has problems when clusters are of differing
- Sizes
- Densities
- Non-globular shapes

◆ K-means has problems when the data contains outliers.

# Limitations of K-means: Differing Sizes



•Original Points

•K-means (3 Clusters)

# Limitations of K-means: Differing Density



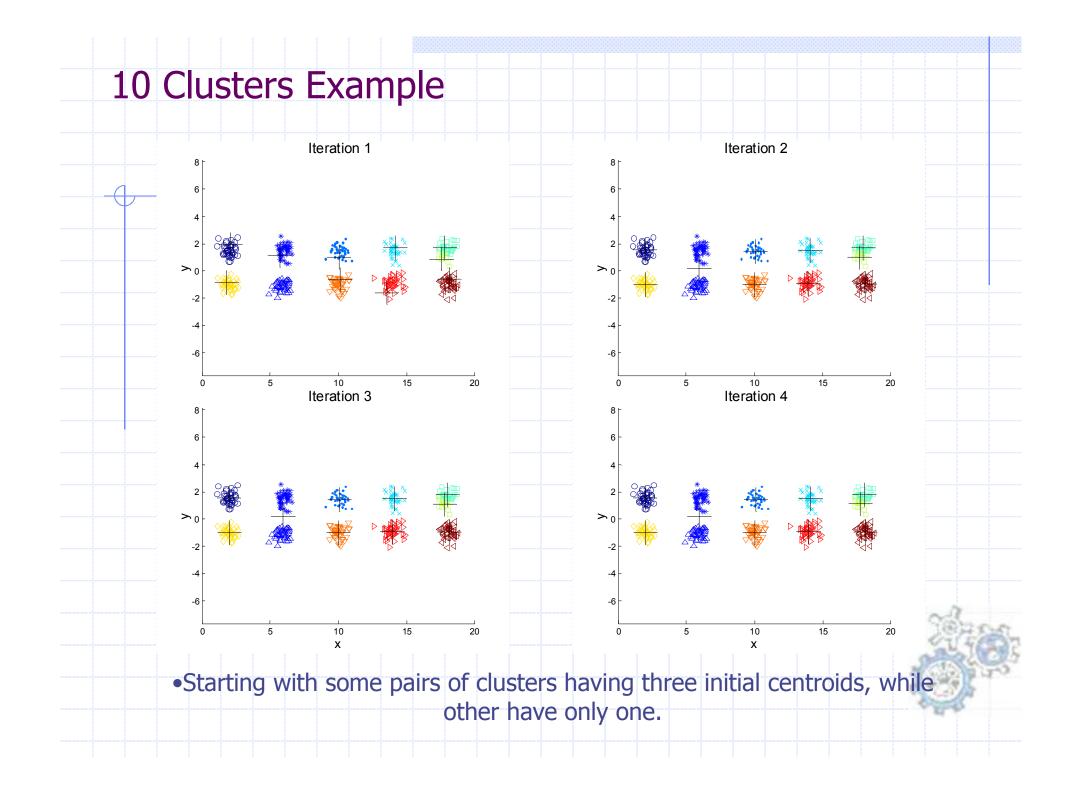- Original Points

- K-means (3 Clusters)

# Limitations of K-means: Non-globular Shapes



•Original Points

•K-means (2 Clusters)

# Overcoming K-means Limitations



•Original Points

K-means Clusters

•One solution is to use many clusters.

•Find parts of clusters, but need to put together.

# Overcoming K-means Limitations



•Original Points                                    K-means Clusters

# Overcoming K-means Limitations



•Original Points                    K-means Clusters

# 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
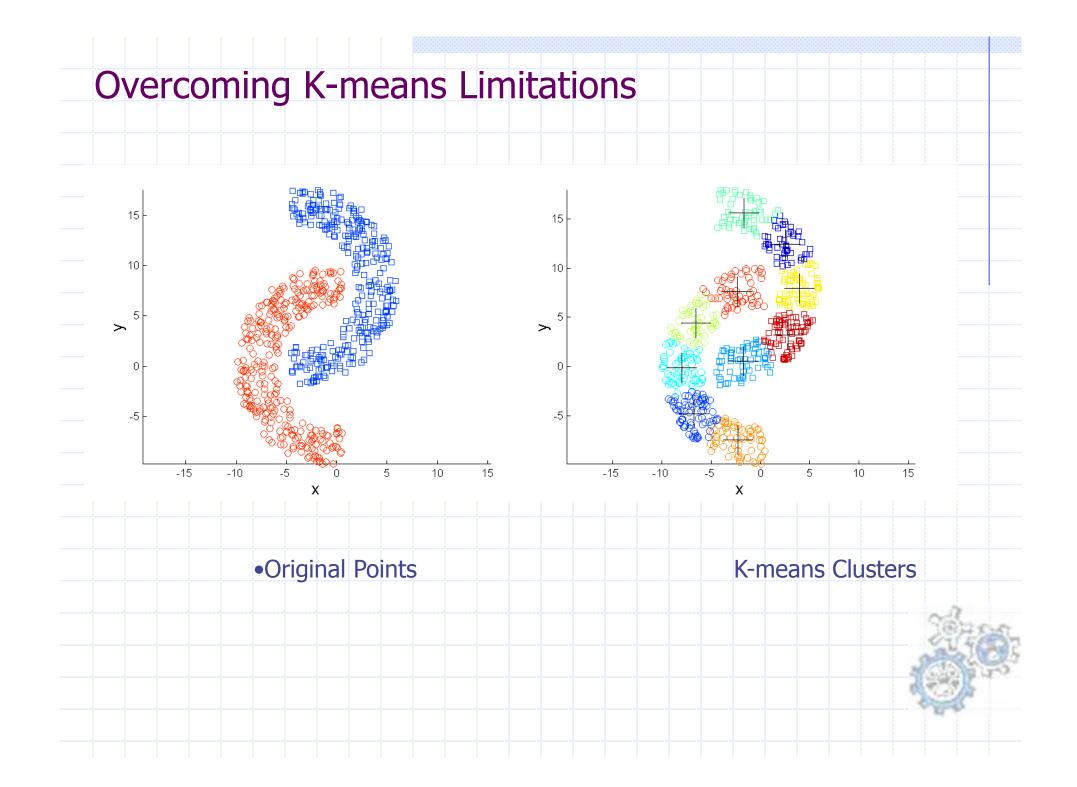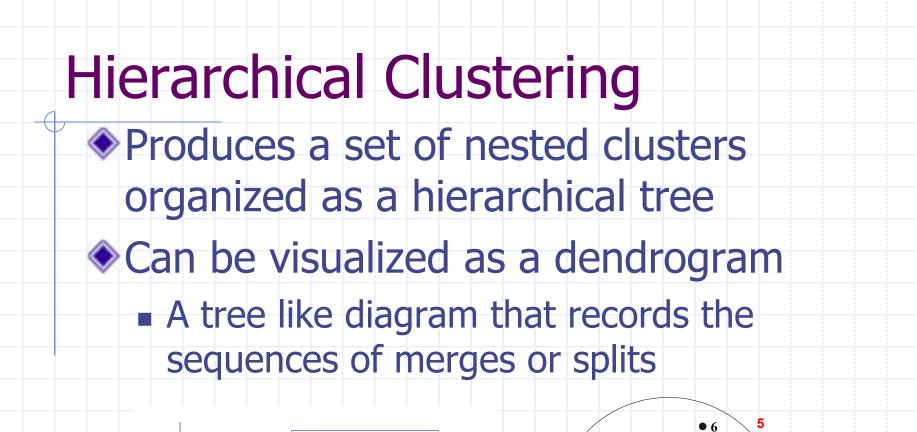  - Any desired number of clusters can be obtained by 'cutting' the dendogram at the proper level

◆ 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

# 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  |     |     |     |     |     |      |
| .    |     |     |     |     |     |      |
| .    |     |     |     |     |     |      |
| .    |     |     |     |     |     |      |

p1    p2    p3    p4    ...    p9    p10    p11    p12

# Intermediate Situation

◆ After some merging steps, we have some clusters

# Intermediate Situation

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

# How to Define Inter-Cluster Similarity

•Similarity?

|  | •p 1 | •p 2 | •p 3 | •p 4 | •p 5 | •. . |
|---|---|---|---|---|---|---|
| •p 1 |  |  |  |  |  | . |
| •p 2 |  |  |  |  |  |  |
| •p 3 |  |  |  |  |  |  |
| •p 4 |  |  |  |  |  |  |
| •p 5 |  |  |  |  |  |  |
| . |  |  |  |  |  |  |
| •. |  |  |  |  |  |  |
| •. |  |  |  |  |  |  |

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

•Proximity Matrix

# How to Define Inter-Cluster Similarity

|        | •p 1 | •p 2 | •p 3 | •p 4 | •p 5 | •.. |
|--------|------|------|------|------|------|-----|
| •p 1   |      |      |      |      |      | .   |
| •p 2   |      |      |      |      |      |     |
| •p 3   |      |      |      |      |      |     |
| •p 4   |      |      |      |      |      |     |
| •p 5   |      |      |      |      |      |     |
| •.     |      |      |      |      |      |     |
| •.     |      |      |      |      |      |     |
| •.     |      |      |      |      |      |     |

•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

|  | •p 1 | •p 2 | •p 3 | •p 4 | •p 5 | •. . |
|---|---|---|---|---|---|---|
| •p 1 |  |  |  |  |  |  |
| •p 2 |  |  |  |  |  |  |
| •p 3 |  |  |  |  |  |  |
| •p 4 |  |  |  |  |  |  |
| •p 5 |  |  |  |  |  |  |
| •. . |  |  |  |  |  |  |

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

•Proximity Matrix

# How to Define Inter-Cluster Similarity



|  | •p 1 | •p 2 | •p 3 | •p 4 | •p 5 | •. . . |
|---|---|---|---|---|---|---|
| •p 1 |  |  |  |  |  |  |
| •p 2 |  |  |  |  |  |  |
| •p 3 |  |  |  |  |  |  |
| •p 4 |  |  |  |  |  |  |
| •p 5 |  |  |  |  |  |  |
| •. |  |  |  |  |  |  |
| •. |  |  |  |  |  |  |
| •. |  |  |  |  |  |  |

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

•Proximity Matrix

# How to Define Inter-Cluster Similarity

|  | •p 1 | •p 2 | •p 3 | •p 4 | •p 5 | •.. . |
|---|---|---|---|---|---|---|
| •p 1 |  |  |  |  |  |  |
| •p 2 |  |  |  |  |  |  |
| •p 3 |  |  |  |  |  |  |
| •p 4 |  |  |  |  |  |  |
| •p 5 |  |  |  |  |  |  |
| •. . |  |  |  |  |  |  |

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

•Proximity Matrix

# 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
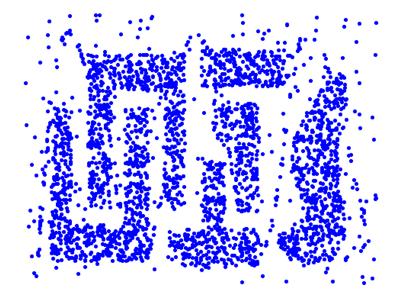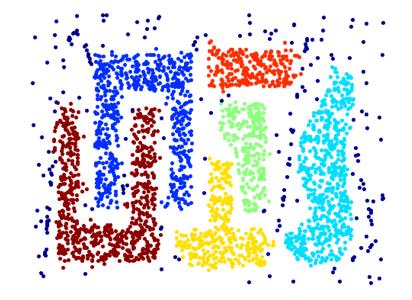  - Breaking large clusters

# 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.
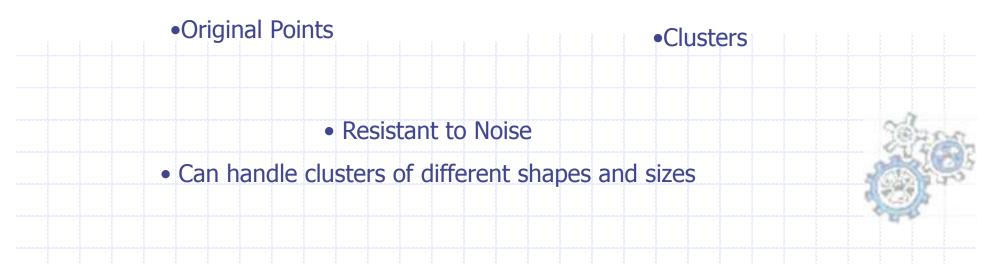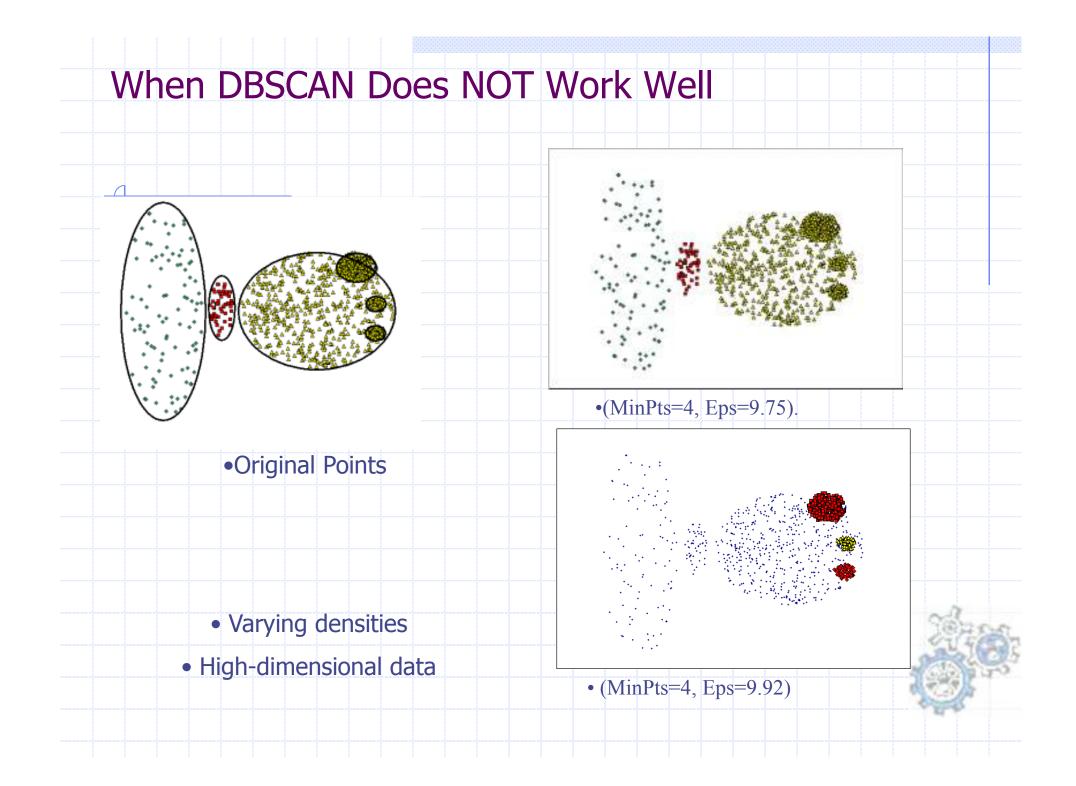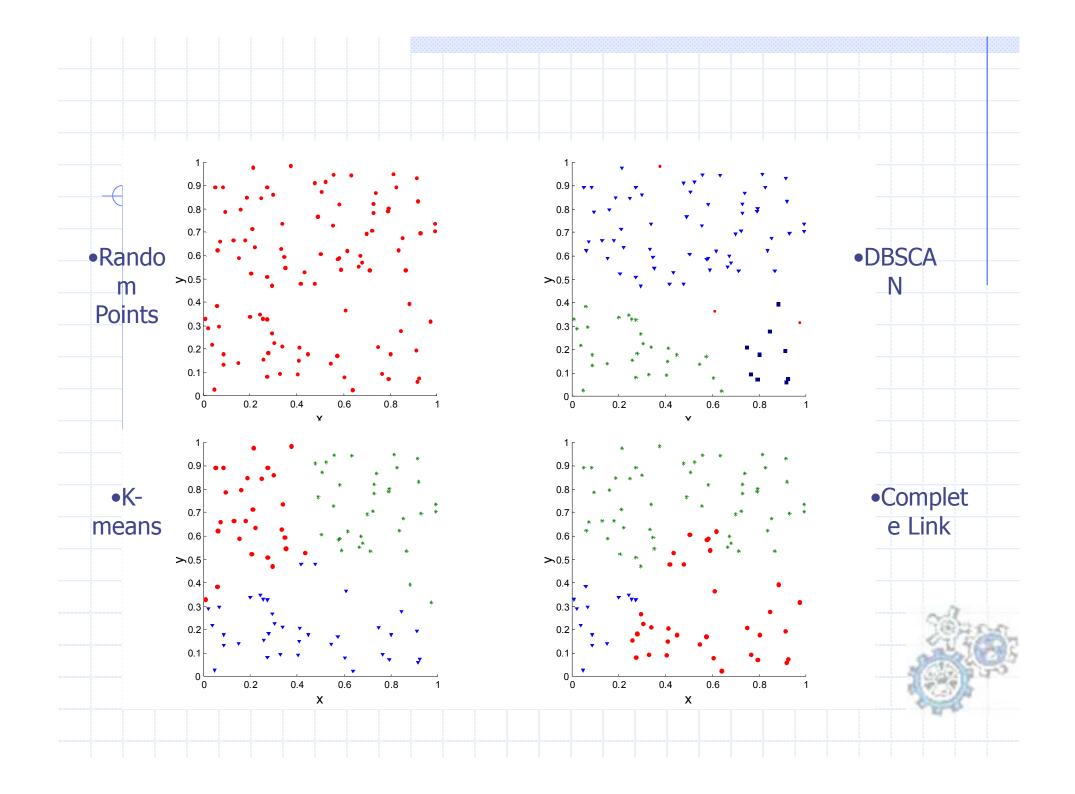
# DBSCAN: Core, Border, and Noise Points

# DBSCAN: Core, Border and Noise Points



- Original Points

- Point types: core, border and noise

- 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).

• Varying densities

• High-dimensional data



• (MinPts=4, Eps=9.92)

# 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

•Rando
m
Points

•DBSCA
N

•K-
means

•Complet
e Link

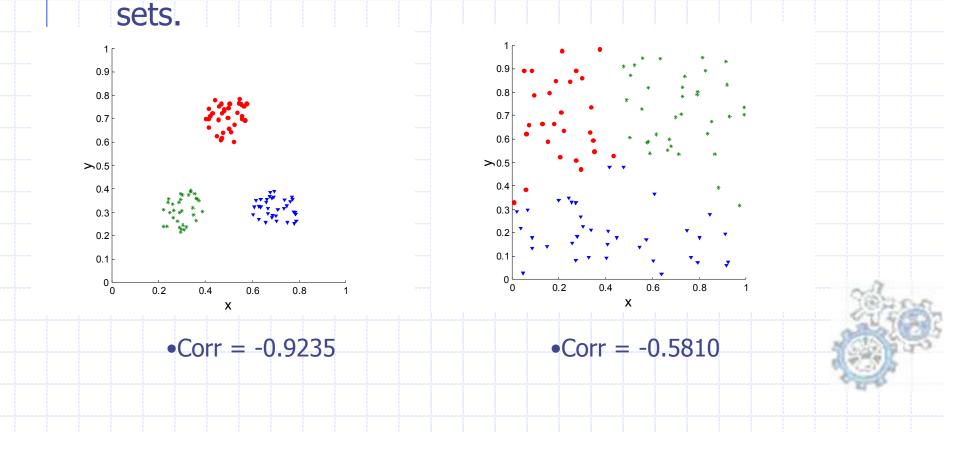# Measuring Cluster Validity Via Correlation

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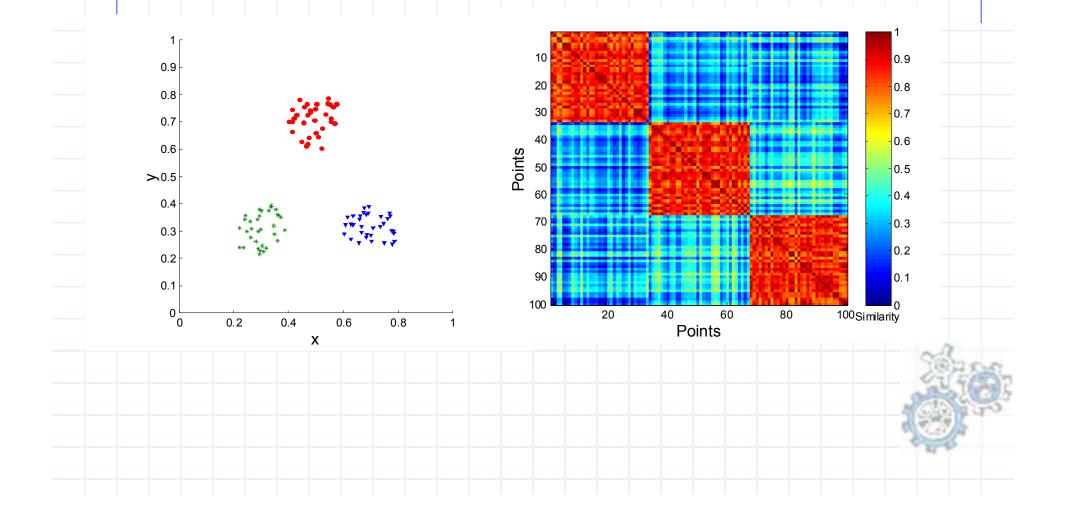- Not a good measure for some density or contiguity based clusters.

# Measuring Cluster Validity Via Correlation

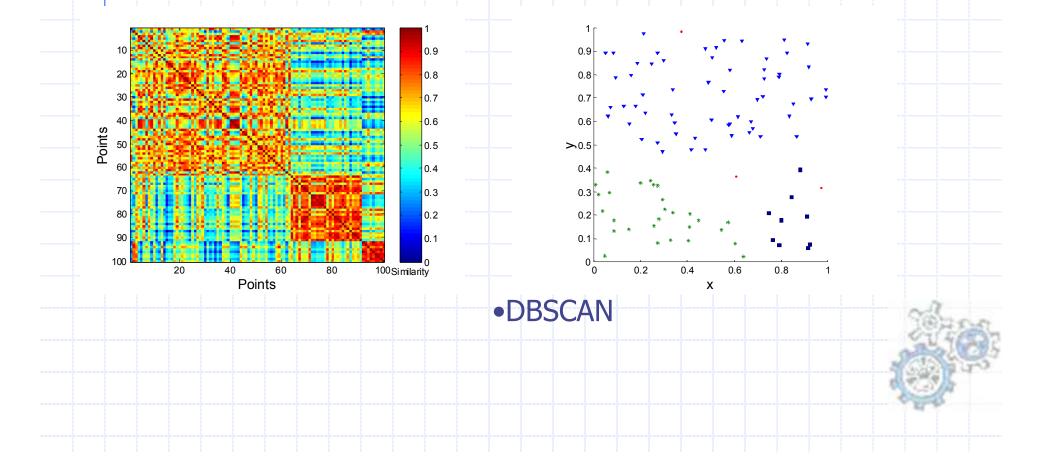◆ Correlation of incidence and proximity matrices for the K-means clusterings of the following two data sets.



• Corr = -0.9235

• Corr = -0.5810

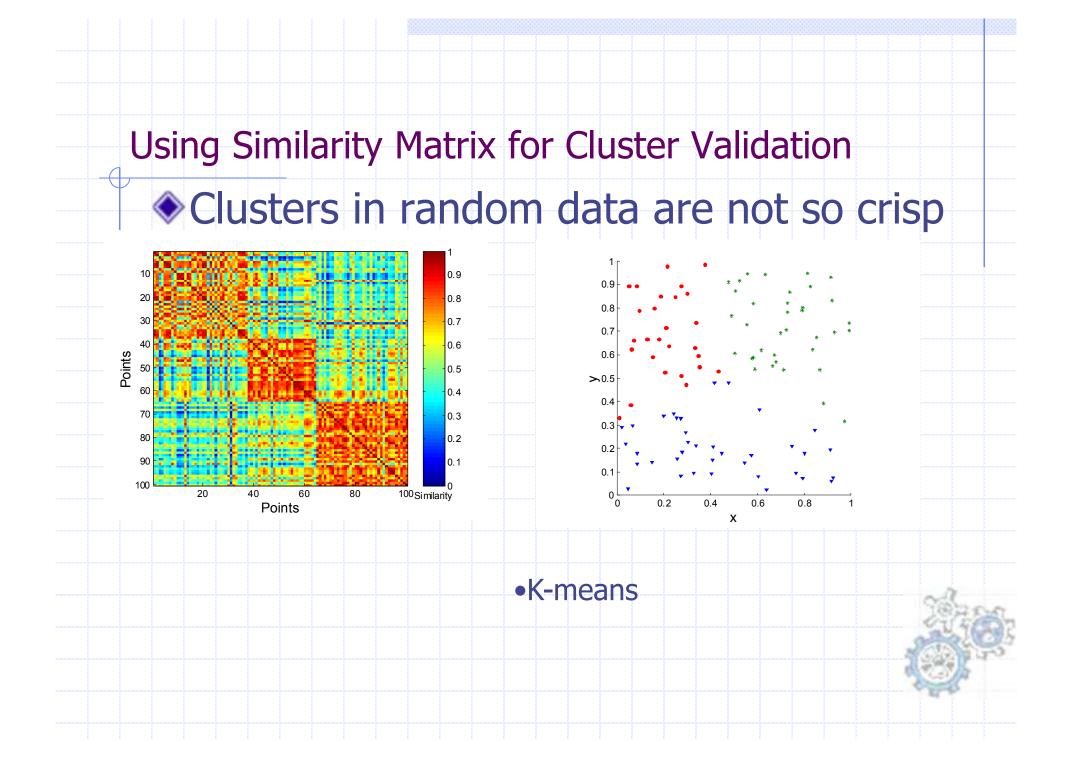# Using Similarity Matrix for Cluster Validation

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

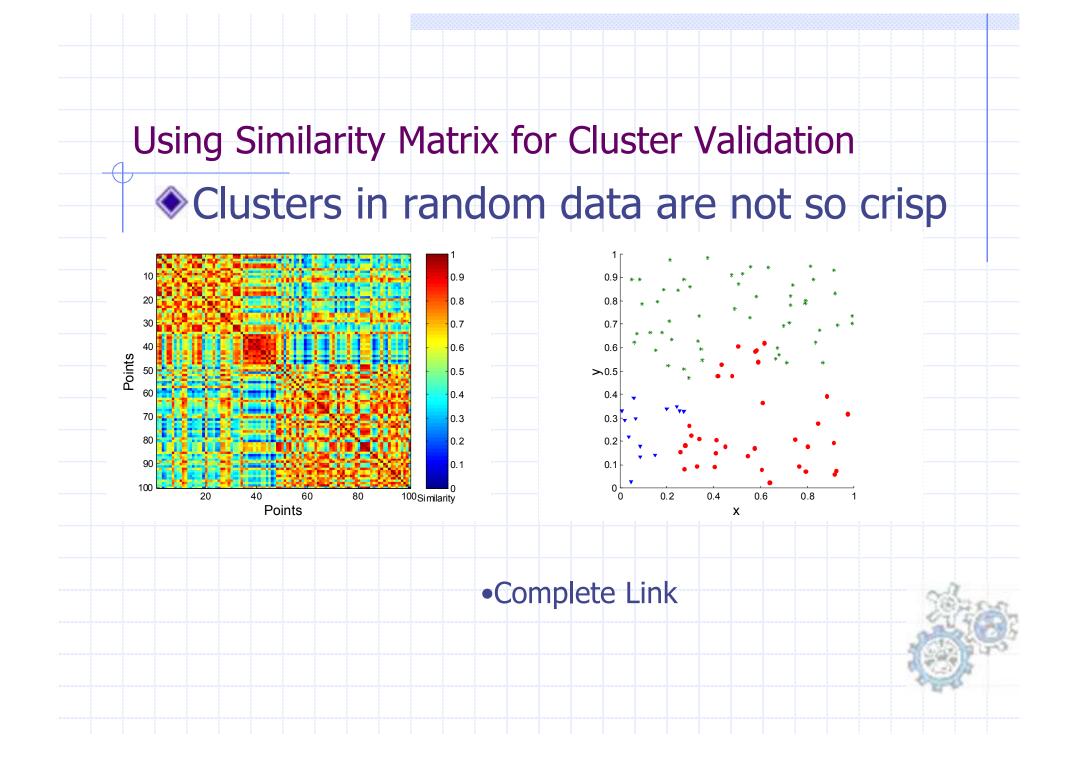# Using Similarity Matrix for Cluster Validation

◆ 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