# Clustering
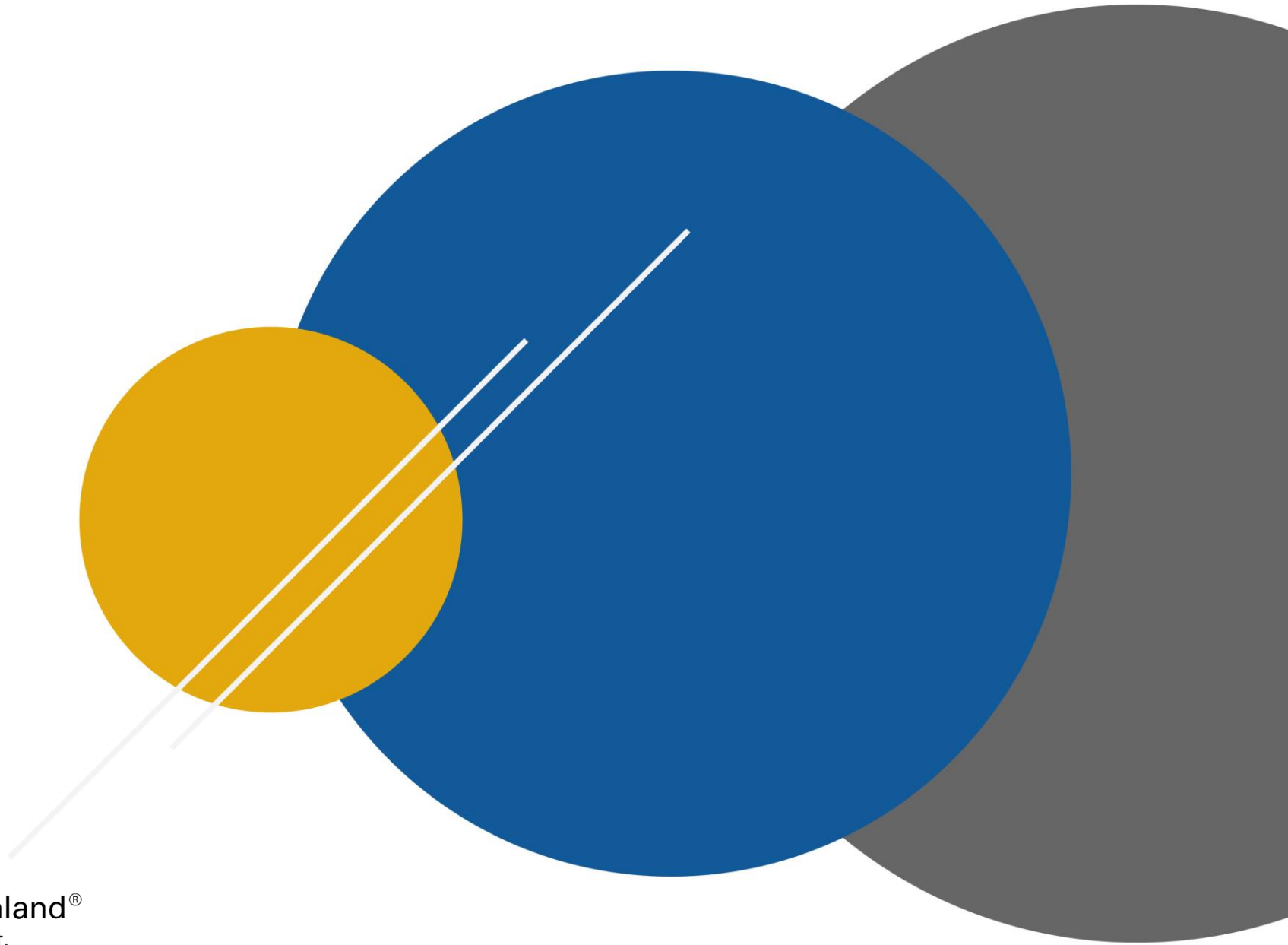Model

Telkom University
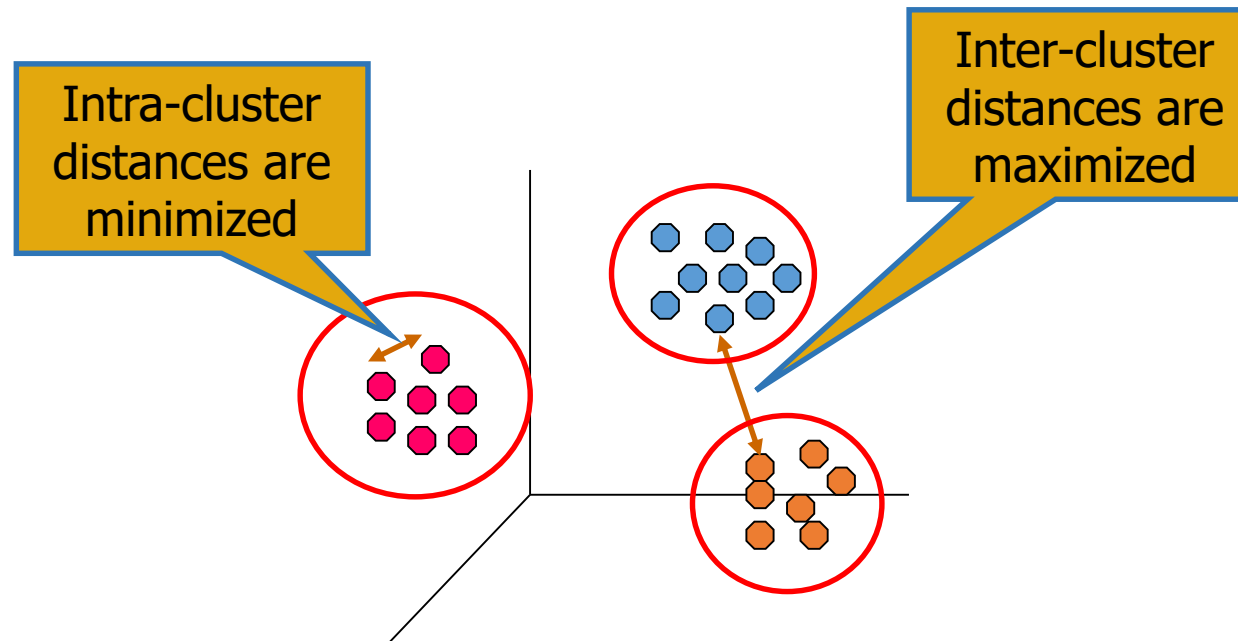
TÜVRheinland®
Precisely Right.

# Agenda

- **Introduction to Clustering**
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Model Evaluation and Selection

# 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

# 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
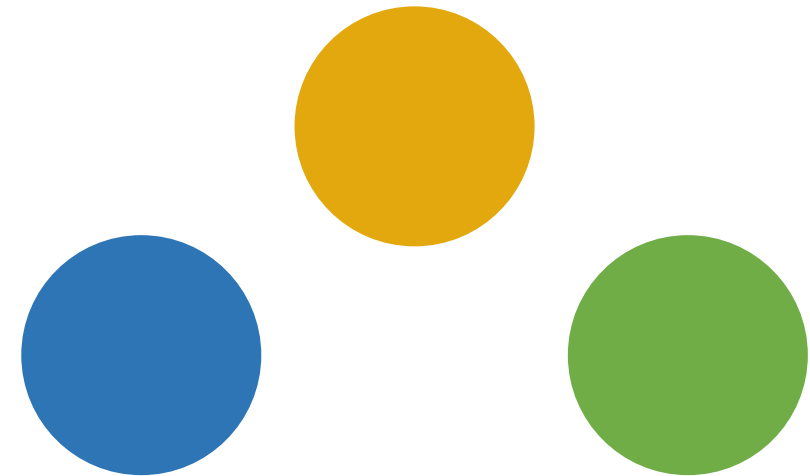  - Some mutual relevance and synergy, but areas are not identical

# Types of Clusters:
## Well-Separated

- 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.
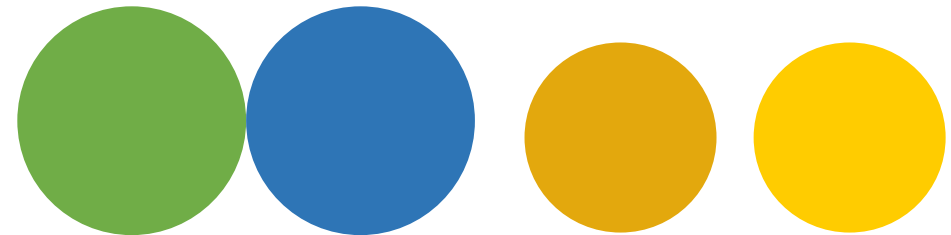
**3 well-separated clusters**

# 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 i.e. the most "representative" point of a cluster
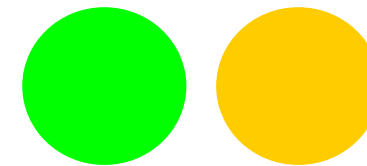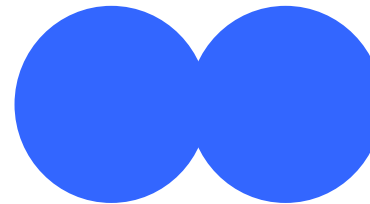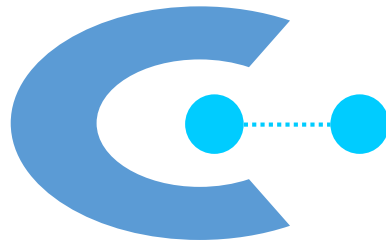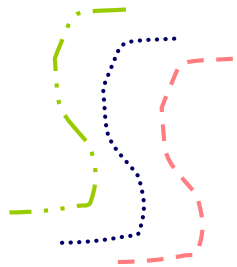
**4 center-based clusters**

# Types of Clusters:
## Contiguity-Based

- 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.
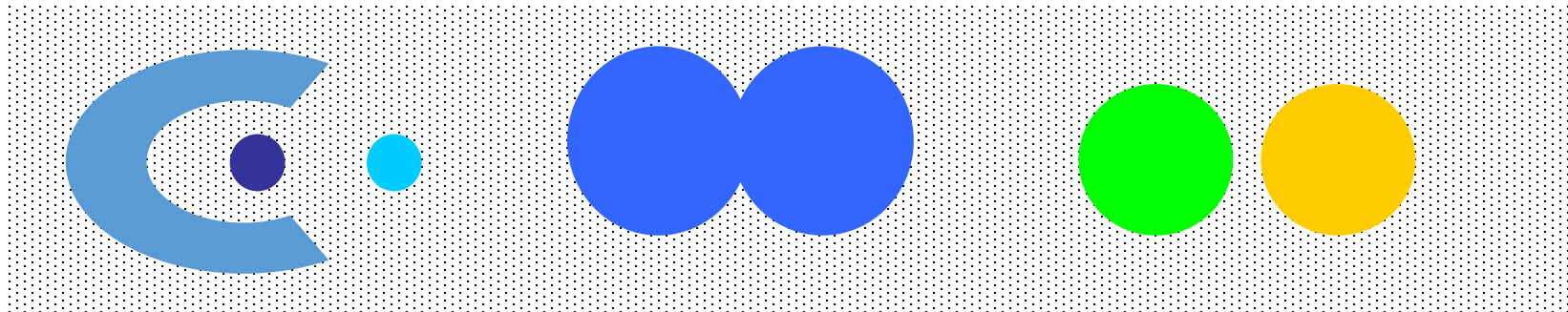
**8 contiguous clusters**

# 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**

# Types of Clusters:
## Objective Function

Clusters Defined by an Objective Function

- Finds clusters that minimize or maximize an objective function.

- Enumerate all possible ways of dividing the points into clusters and evaluate the `goodness' of each potential set of clusters by using the given objective function.  (NP Hard)

- Can have global or local objectives.

  - Hierarchical clustering algorithms typically have local objectives
  - Partitional algorithms typically have global objectives

# Types of Clusters:
## Objective Function

Clusters Defined by an Objective Function

- A variation of the global objective function approach is to fit the data to a parameterized model.
  - Parameters for the model are determined from the data.
  - Mixture models assume that the data is a 'mixture' of a number of statistical distributions.

- Map the clustering problem to a different domain and solve a related problem in that domain
  - Proximity matrix defines a weighted graph, where the nodes are the points being clustered, and the weighted edges represent the proximities between points
  - Clustering is equivalent to breaking the graph into connected components, one for each cluster.
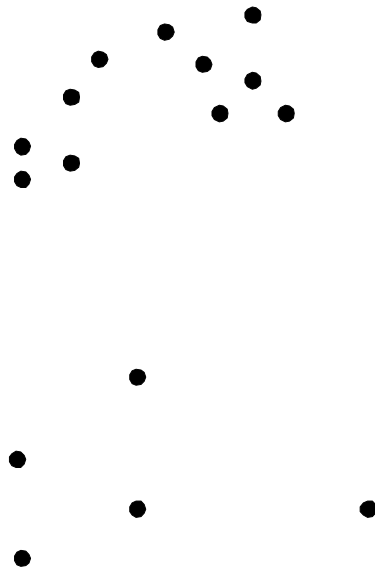  - Want to minimize the edge weight between clusters and maximize the edge weight within clusters
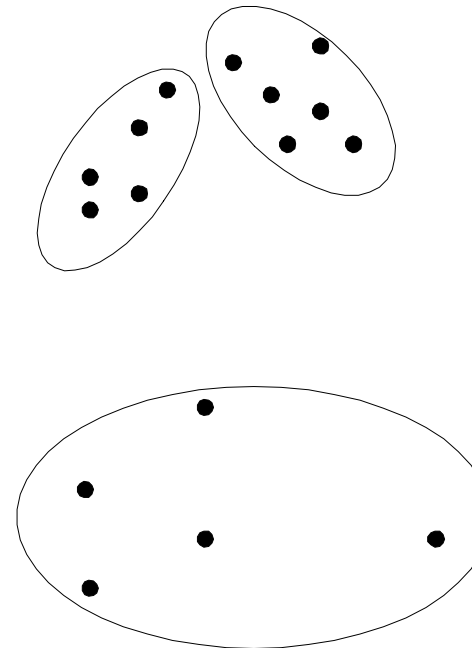
# Agenda

Clustering Model

# Partitional Clustering



**Original Points**           **A Partitional Clustering**

# K-means Algorithm

- 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

---

**Algorithm 1** Basic K-means Algorithm.

---

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.
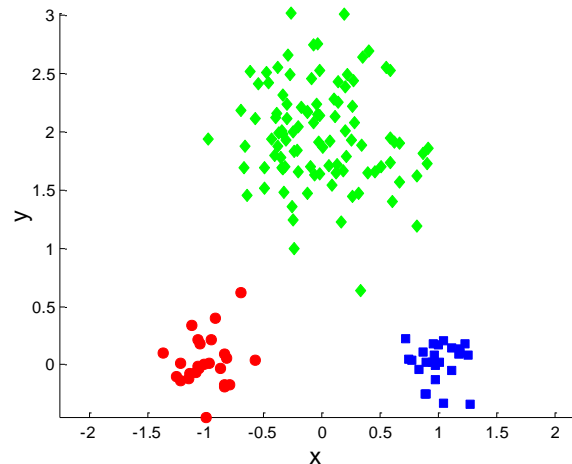
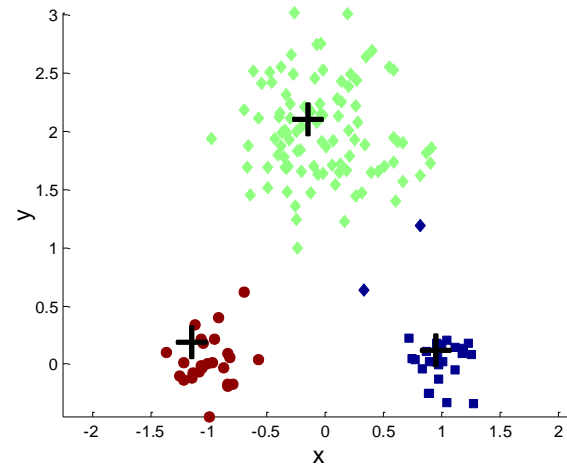5: **until** The centroids don't change

---

# 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 commonly measured by Euclidean distance.
- 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
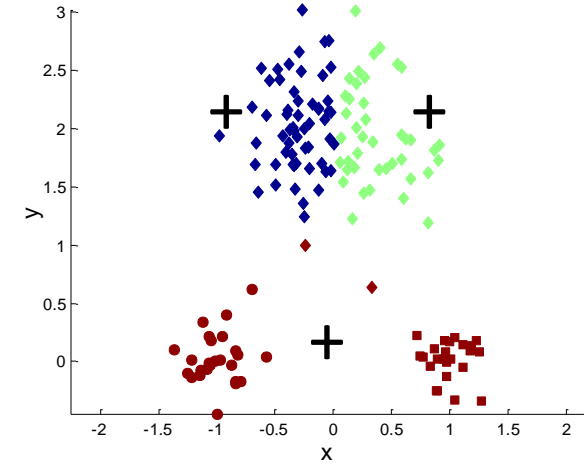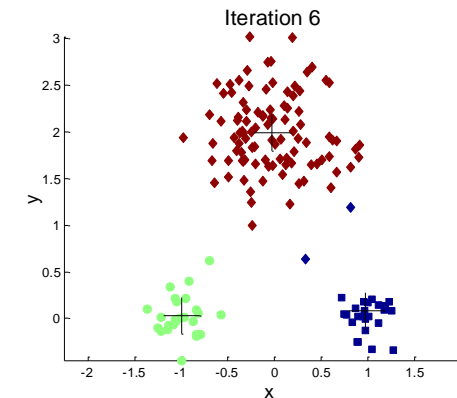
# Two different
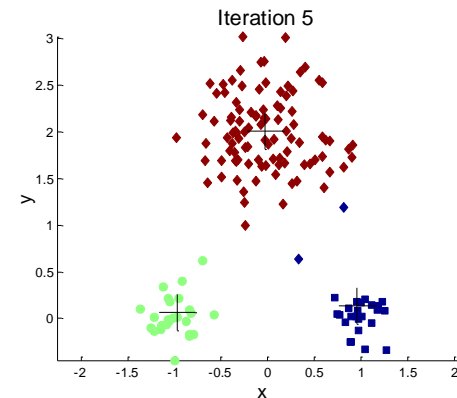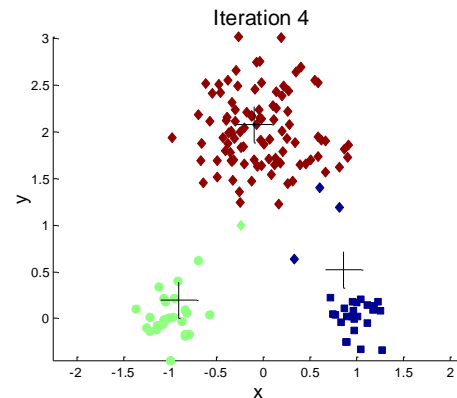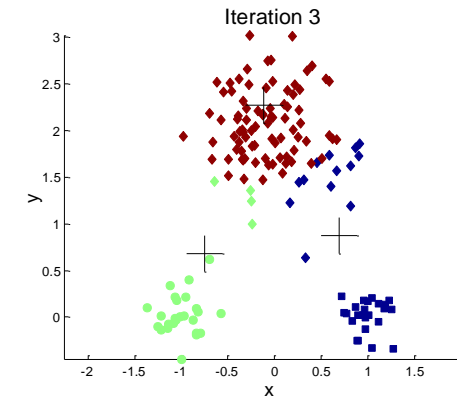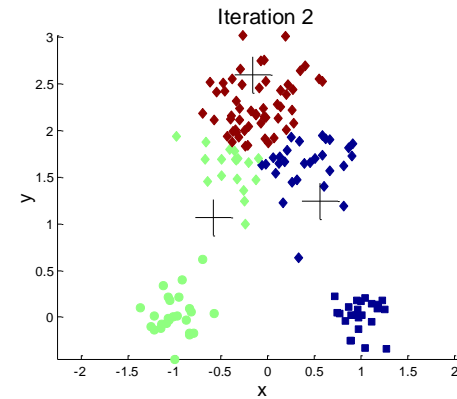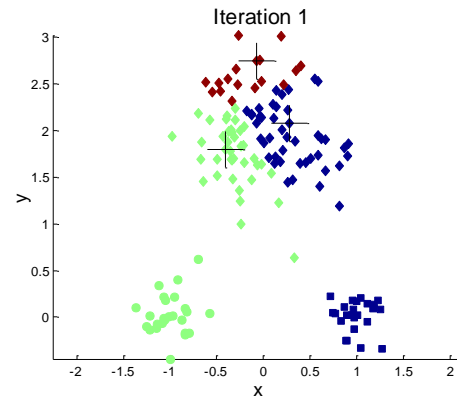## K-means Clusterings



Original Points

Optimal Clustering

Sub-optimal Clustering

# Importance of
# Choosing Initial Centroids

# Evaluating
## K-means Clusters

Most common measure is Sum of Squared Error (SSE)

- For each point, the error is the distance to the nearest cluster

- To get SSE, we square these errors and sum them.

$$SSE = \sum_{i=1}^{K} \sum_{x \in C_i} dist^2(m_i, x)$$

- x is a data point in cluster Ci and mi is the representative point for cluster Ci
  - can show that mi corresponds to the center (mean) of the cluster

- Given two clusters, we can choose the one with the smallest error

- One easy way to reduce SSE is to increase K, the number of clusters
  - A good clustering with smaller K can have a lower SSE than a poor clustering with higher K

# Selecting Optimal
# Number of Cluster

- Elbow Method
  - Compute clustering algorithm (e.g., k-means clustering) for different values of k. For instance, by varying k from 1 to 10 clusters.
  - For each k, calculate the total within-cluster sum of square (wss).
  - Plot the curve of wss according to the number of clusters k.
  - The location of a bend (knee) in the plot is generally considered as an indicator of the appropriate number of clusters.

- Average silhouette method
  - Compute clustering algorithm (e.g., k-means clustering) for different values of k. For instance, by varying k from 1 to 10 clusters.
  - For each k, calculate the average silhouette of observations (avg.sil).
  - Plot the curve of avg.sil according to the number of clusters k.
  - The location of the maximum is considered as the appropriate number of clusters.

# Selecting Optimal
# Number of Cluster

These methods include direct methods and statistical testing methods:

- Direct methods: consists of optimizing a criterion, such as the within cluster sums of squares or the average silhouette. The corresponding methods are named elbow and silhouette methods, respectively.

- Statistical testing methods: consists of comparing evidence against null hypothesis. An example is the gap statistic.

Clustering Model

# Problems with
# Selecting Initial Points

If there are K 'real' clusters then the chance of selecting one centroid from each cluster is small.

- Chance is relatively small when K is large

- If clusters are the same size, n, then

$$P = \frac{\text{number of ways to select one centroid from each cluster}}{\text{number of ways to select } K \text{ centroids}} = \frac{K!n^K}{(Kn)^K} = \frac{K!}{K^K}$$

For example, if K = 10, then probability = 10!/1010 = 0.00036

- Sometimes the initial centroids will readjust themselves in 'right' way, and sometimes they don't

- Consider an example of five pairs of clusters

# Solutions to
# Initial Centroids Problem

- Multiple runs
  - Helps, but probability is not on your side
- Sample and use hierarchical clustering to determine initial centroids
- Select more than k initial centroids and then select among these initial centroids
  - Select most widely separated
- Post-processing
- Bisecting K-means
  - Not as susceptible to initialization issues

# Pre-processing and
# Post-processing
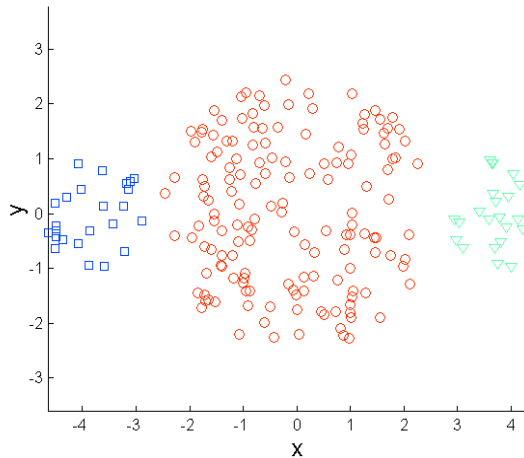
Pre-processing

- Normalize the data
- Eliminate outliers

Post-processing

- Eliminate small clusters that may represent outliers
- Split 'loose' clusters, i.e., clusters with relatively high SSE
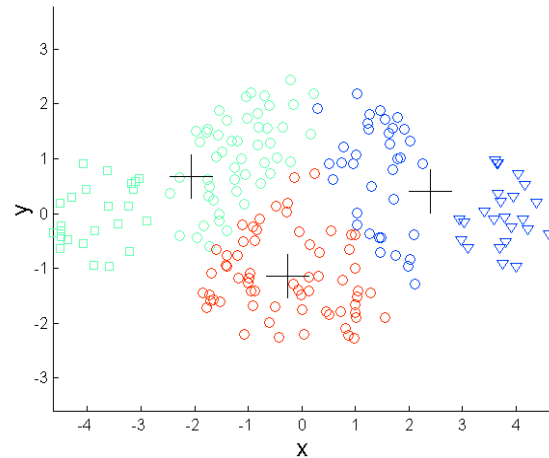- Merge clusters that are 'close' and that have relatively low SSE

# Limitations & Overcoming of K-means:
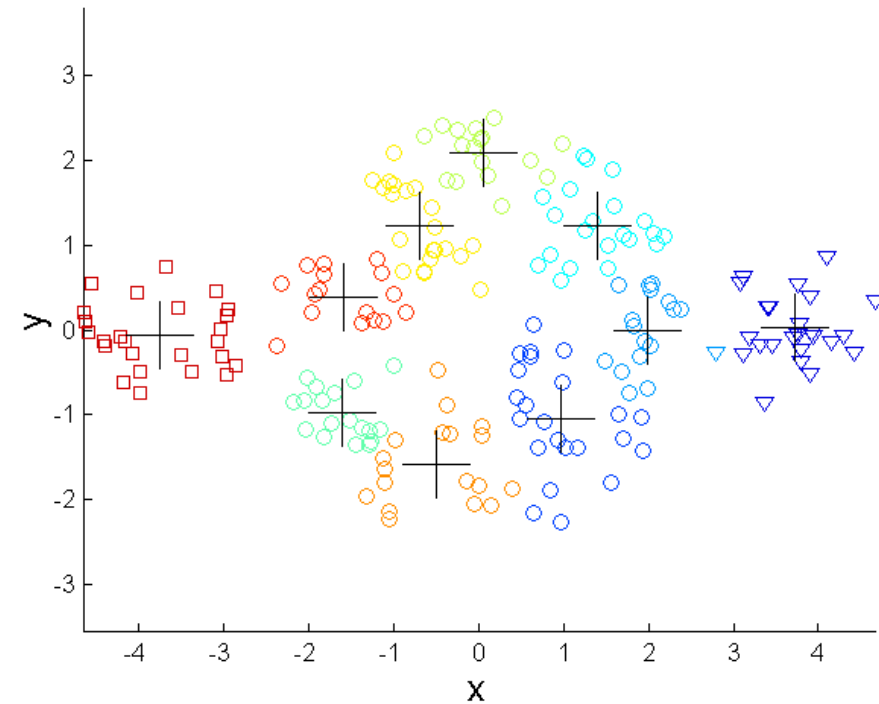## Differing Sizes

One solution is to use many clusters.
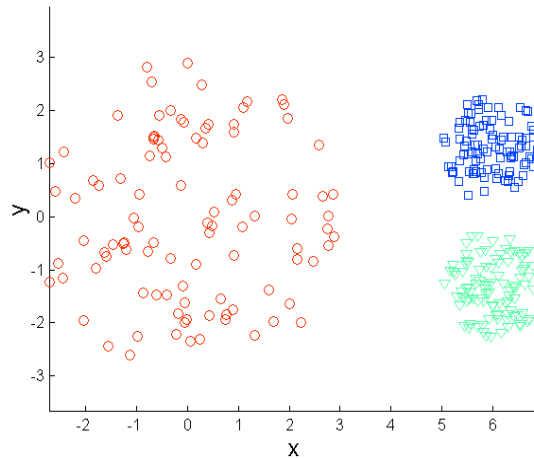Find parts of clusters, but need to put together.
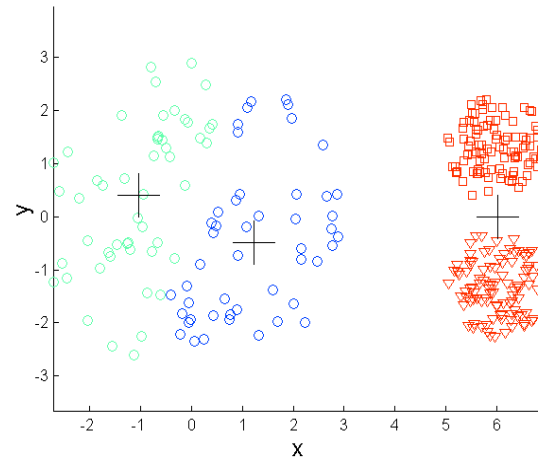
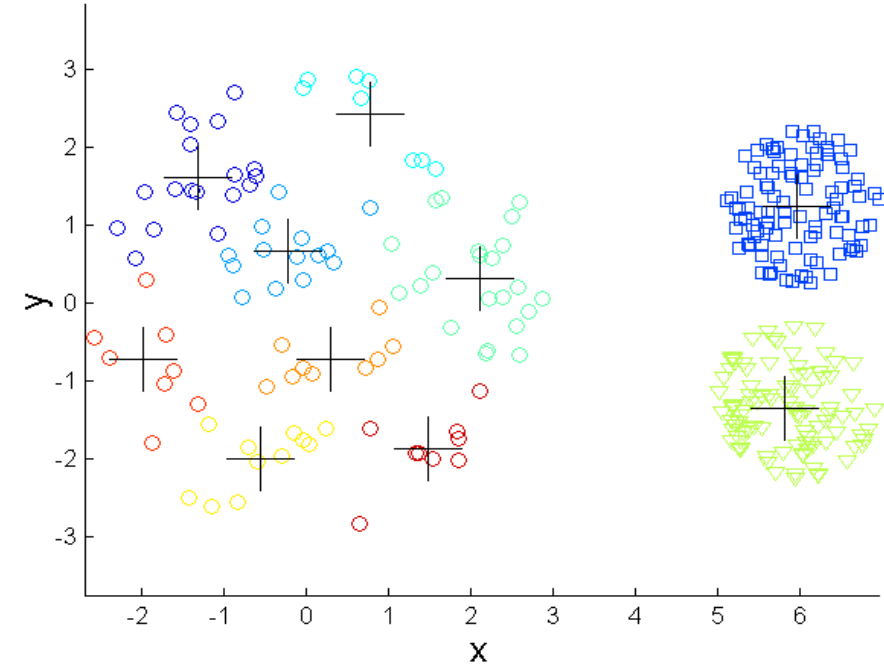

Original Points

K-means (3 Clusters)

Overcoming K-means

# Limitations & Overcoming of K-means:
## Differing Density
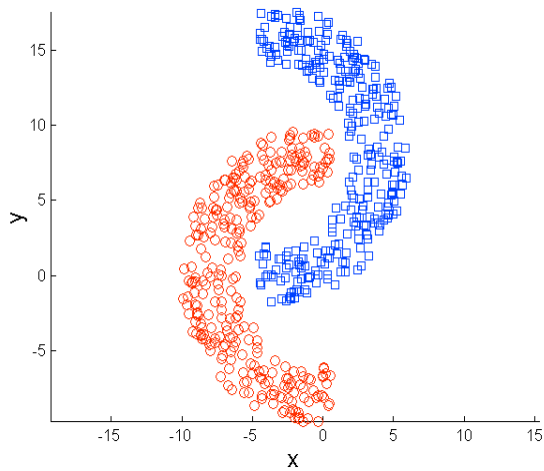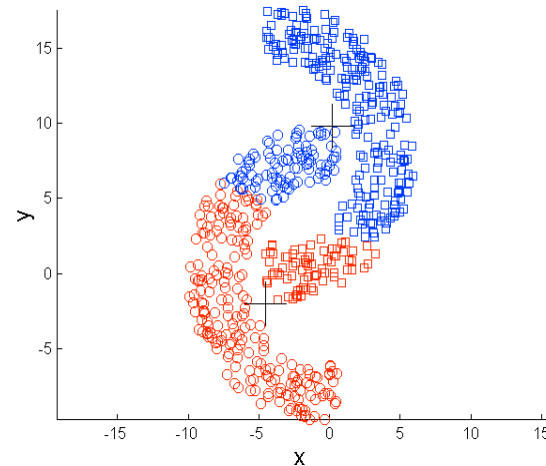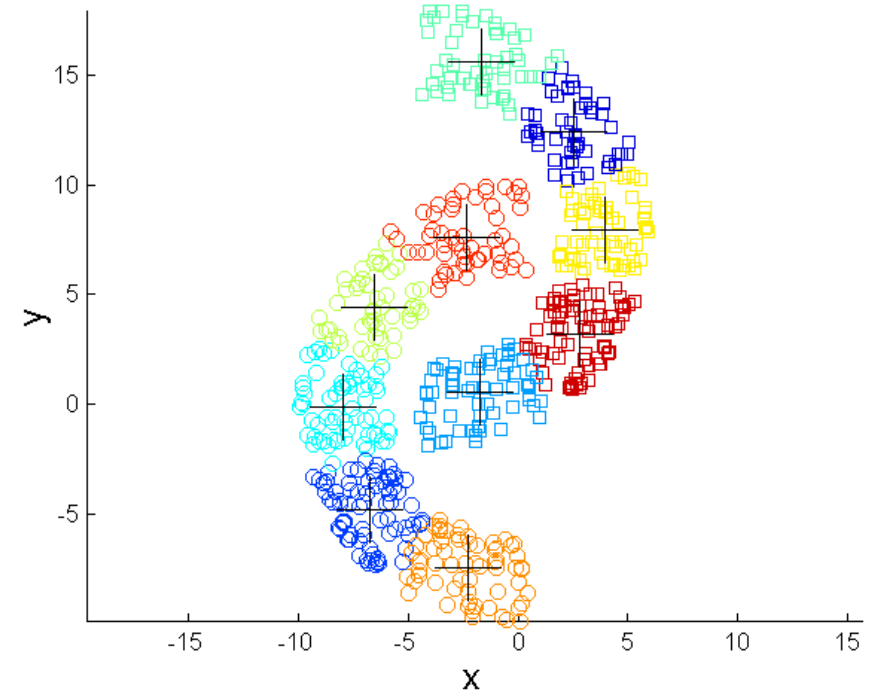


Original Points

K-means (3 Clusters)

Overcoming K-means

# Limitations & Overcoming of K-means:
# Non-globular Shapes
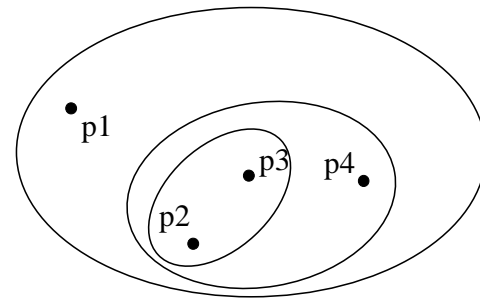


Original Points
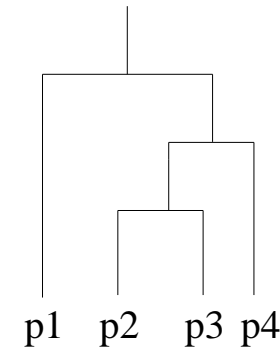
K-means (2 Clusters)

Overcoming K-means
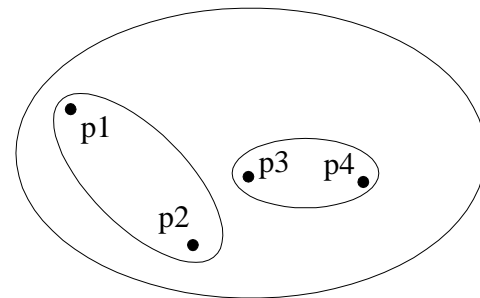
# Agenda

Clustering Model
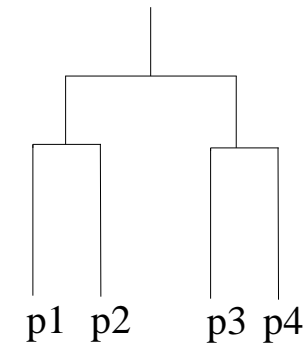
# Hierarchical Clustering



Traditional Hierarchical Clustering

Traditional dendogram

Non-traditional Hierarchical Clustering

Non-traditional dendogram

# Hierarchical Clustering

- Produces a set of nested clusters organized as a hierarchical tree
  - Can be visualized as a dendogram
  - 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



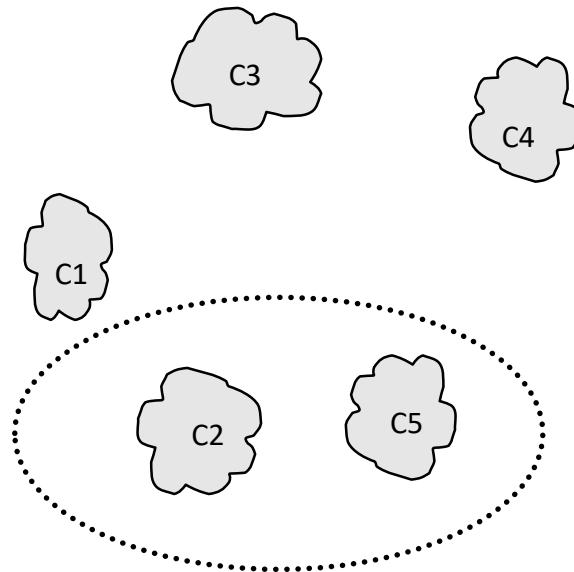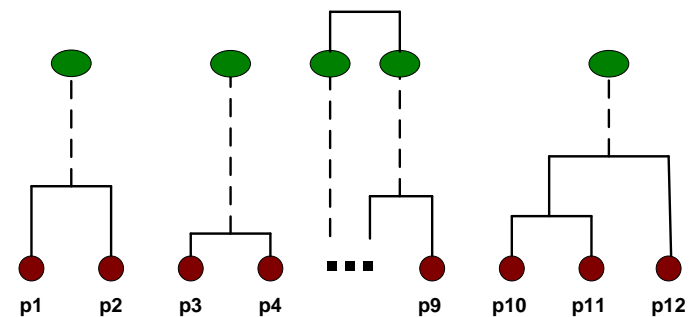| | C1 | C2 | C3 | C4 | C5 |
|---|---|---|---|---|---|
| C1 | | | | | |
| C2 | | | | | |
| C3 | | | | | |
| C4 | | | | | |
| C5 | | | | | |

Proximity Matrix

# Intermediate Situation

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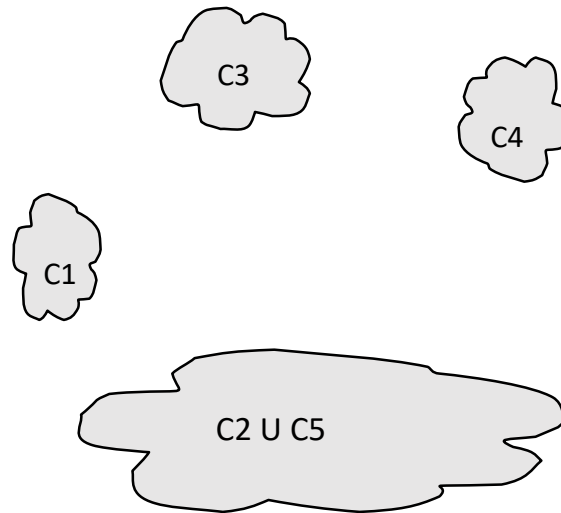


Proximity Matrix

# After Merging

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

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

Proximity Matrix



C3

C4

C1

C2 U C5

p1  p2  p3  p4  p9  p10  p11  p12

# How to Define
## Inter-Cluster Similarity

- MIN

- MAX

- Group Average

- Distance Between Centroids

- Other methods driven by an objective function

  - Ward's Method uses squared error

Similarity?

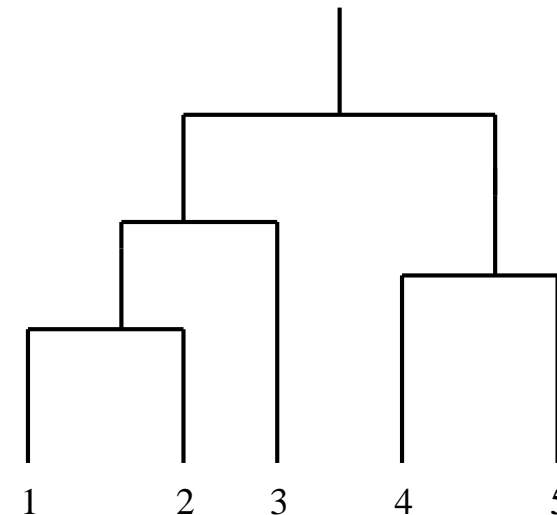|    | p1 | p2 | p3 | p4 | p5 | . . . |
|----|----|----|----|----|----|-------|
| p1 |    |    |    |    |    |       |
| p2 |    |    |    |    |    |       |
| p3 |    |    |    |    |    |       |
| p4 |    |    |    |    |    |       |
| p5 |    |    |    |    |    |       |
| .  |    |    |    |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

# 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.

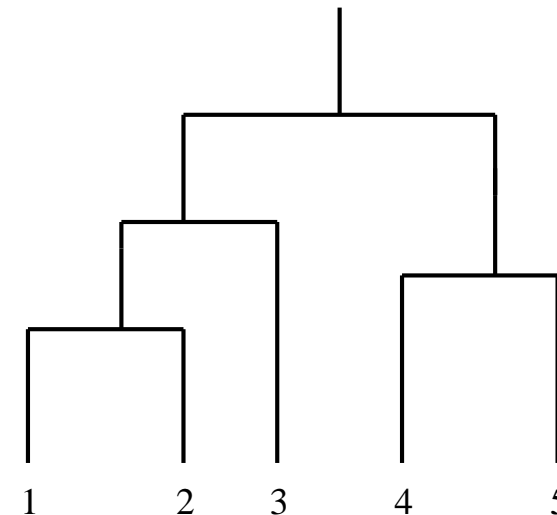|    | I1 | I2 | I3 | I4 | I5 |
|----|------|------|------|------|------|
| I1 | 1.00 | 0.90 | 0.10 | 0.65 | 0.20 |
| I2 | 0.90 | 1.00 | 0.70 | 0.60 | 0.50 |
| I3 | 0.10 | 0.70 | 1.00 | 0.40 | 0.30 |
| I4 | 0.65 | 0.60 | 0.40 | 1.00 | 0.80 |
| I5 | 0.20 | 0.50 | 0.30 | 0.80 | 1.00 |



*Proximity matrix is based on Correlation

# 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

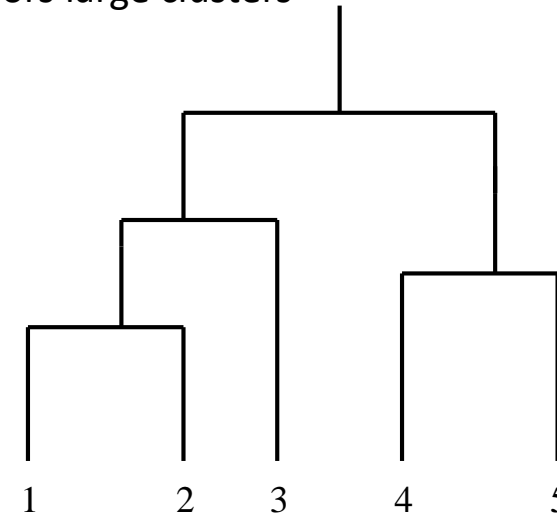|    | I1 | I2 | I3 | I4 | I5 |
|----|----|----|----|----|----|
| I1 | 1.00 | 0.90 | 0.10 | 0.65 | 0.20 |
| I2 | 0.90 | 1.00 | 0.70 | 0.60 | 0.50 |
| I3 | 0.10 | 0.70 | 1.00 | 0.40 | 0.30 |
| I4 | 0.65 | 0.60 | 0.40 | 1.00 | 0.80 |
| I5 | 0.20 | 0.50 | 0.30 | 0.80 | 1.00 |

# Cluster Similarity:
## Group Average

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

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

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

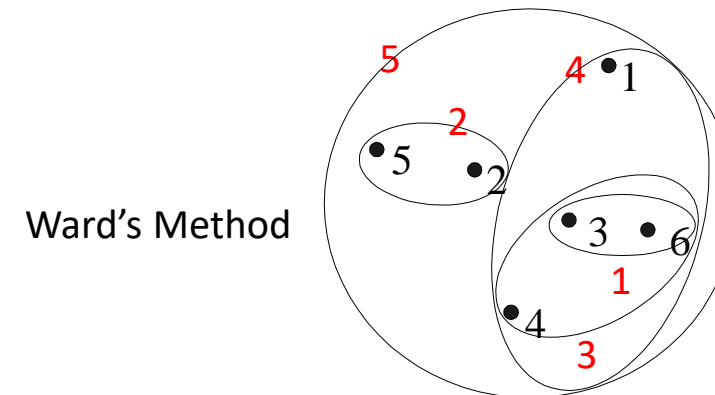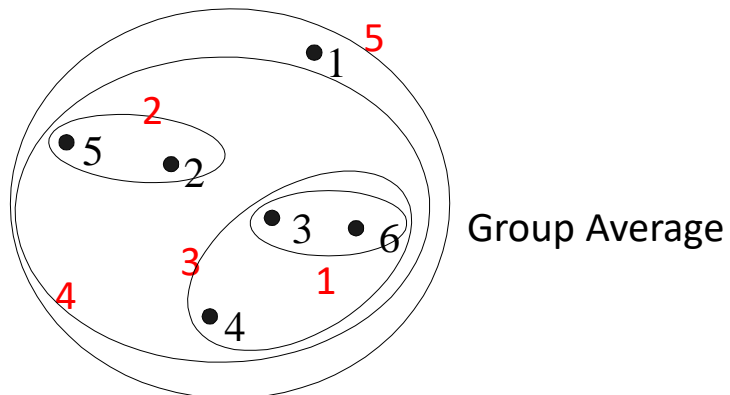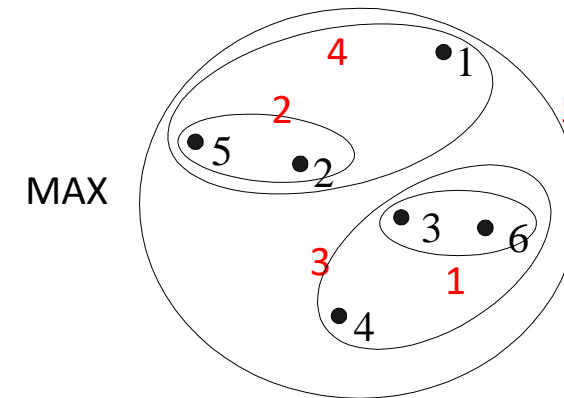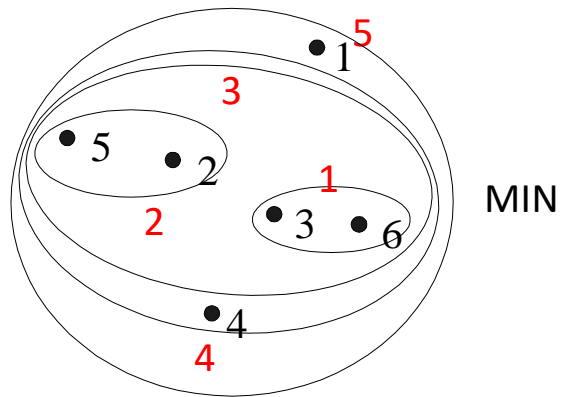|    | I1   | I2   | I3   | I4   | I5   |
|----|------|------|------|------|------|
| I1 | 1.00 | 0.90 | 0.10 | 0.65 | 0.20 |
| I2 | 0.90 | 1.00 | 0.70 | 0.60 | 0.50 |
| I3 | 0.10 | 0.70 | 1.00 | 0.40 | 0.30 |
| I4 | 0.65 | 0.60 | 0.40 | 1.00 | 0.80 |
| I5 | 0.20 | 0.50 | 0.30 | 0.80 | 1.00 |

# 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:
## Comparison

**STRENGTH**

- MIN :

  - Can handle non-elliptical shapes

- MAX :

  - Less susceptible to noise and outliers

- GROUP AVG :

  - Less susceptible to noise and outliers

**LIMITATION**

- MIN:

  - Sensitive to noise and outliers

- MAX :

  - Tends to break large clusters

  - Biased towards globular clusters

- GROUP AVG :

  - Biased towards globular clusters

# 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
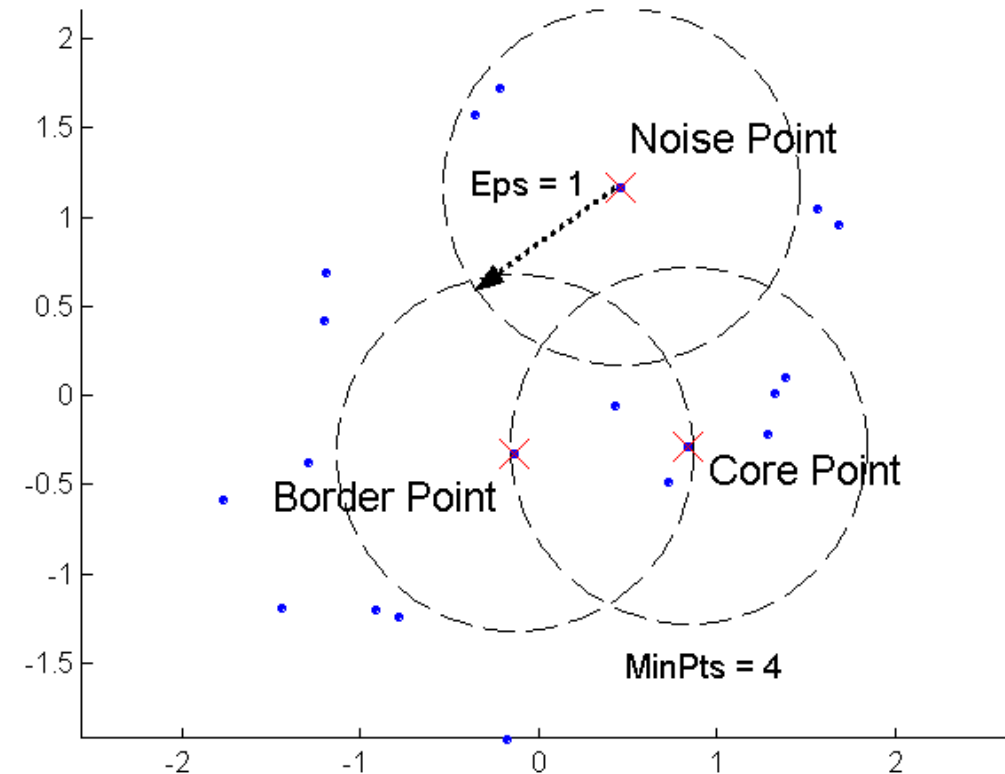
  - Breaking large clusters

# Agenda

Introduction to Clustering

Partitioning Methods

Hierarchical Methods

**Density-Based Methods**

Model Evaluation and Selection

# 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
# 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 $current\_cluster\_label$

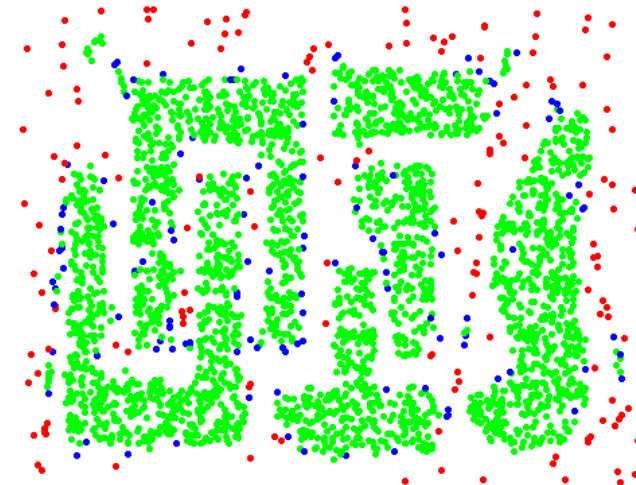        **end if**

    **end for**

**end for**

# DBSCAN:
## Core, Border and Noise Points
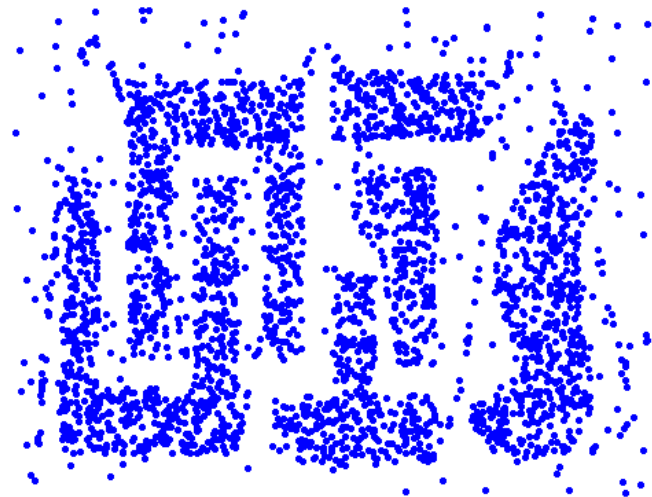


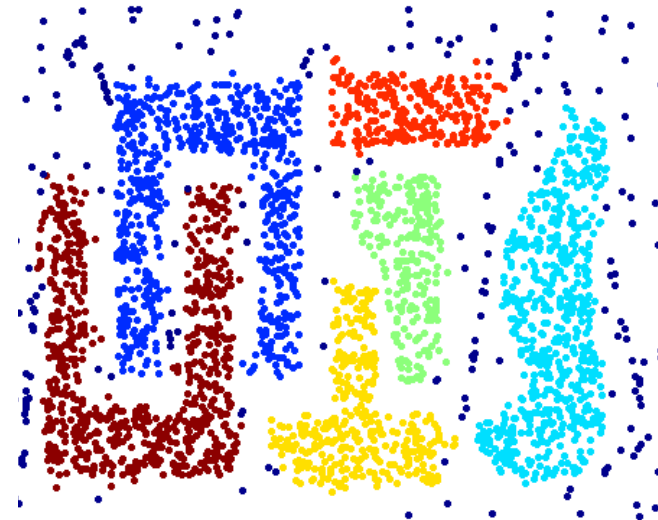Original Points

Point types: core, border and noise

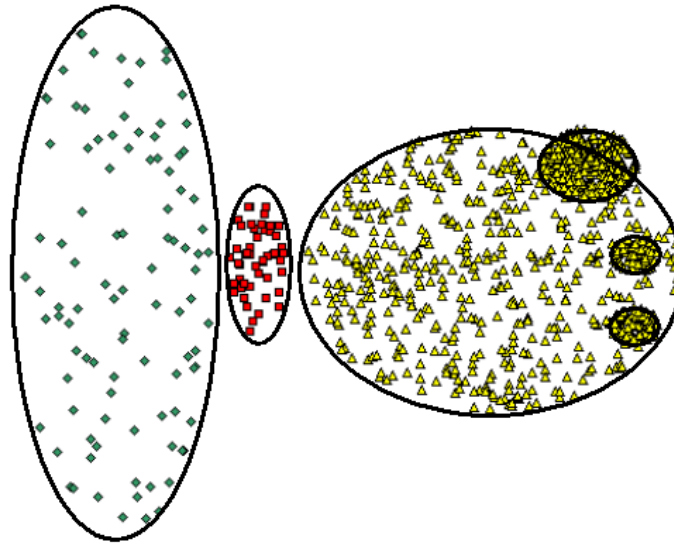Eps = 10, MinPts = 4

# When DBSCAN
Works Well



Original Points



Clusters

Resistant to Noise

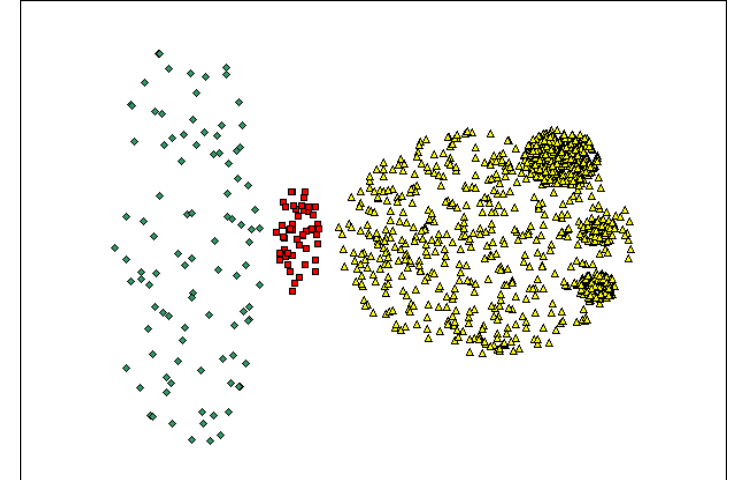Can handle clusters of different shapes and sizes
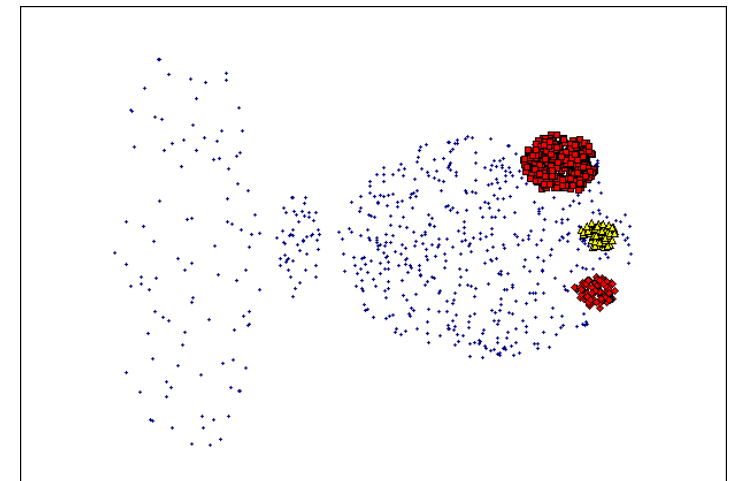
# When DBSCAN
## Does NOT Works Well



Original Points

- Varying densities
- High-dimensional data



(MinPts=4, Eps=9.75).



(MinPts=4, Eps=9.92)

Clustering Model

# Agenda

- Introduction to Clustering
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- **Model Evaluation and Selection**

# 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

# Different Aspects of
# Cluster Validation

- Determining the clustering tendency of a set of data, i.e., distinguishing whether non-random structure actually exists in the data.

- Comparing the results of a cluster analysis to externally known results, e.g., to externally given class labels.

- Evaluating how well the results of a cluster analysis fit the data without reference to external information.
  - Use only the data

- Comparing the results of two different sets of cluster analyses to determine which is better.

- Determining the 'correct' number of clusters.

For 2, 3, and 4, we can further distinguish whether we want to evaluate the entire clustering or just individual clusters.

# 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.
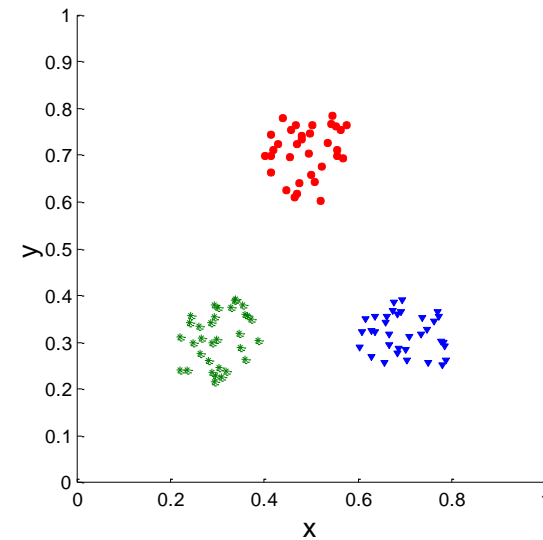
# Validity Measurement
## Via Correlation

- Two matrices
  - Proximity Matrix – distance between any pair of points
  - "Incidence" Matrix – association pair of points
    - 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.

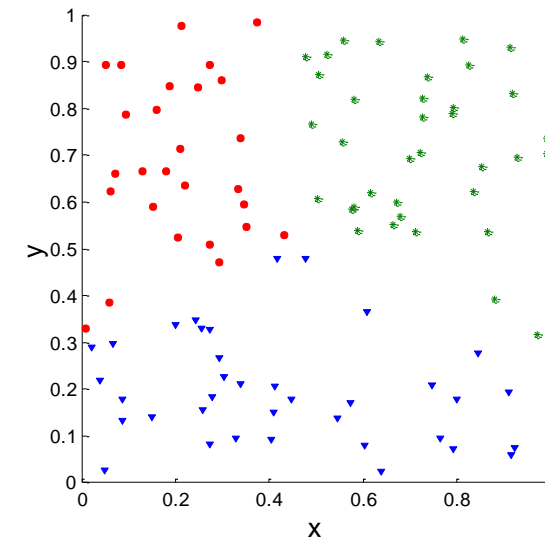- Not a good measure for some density or contiguity based clusters.

# Validity Measurement
## Via Correlation

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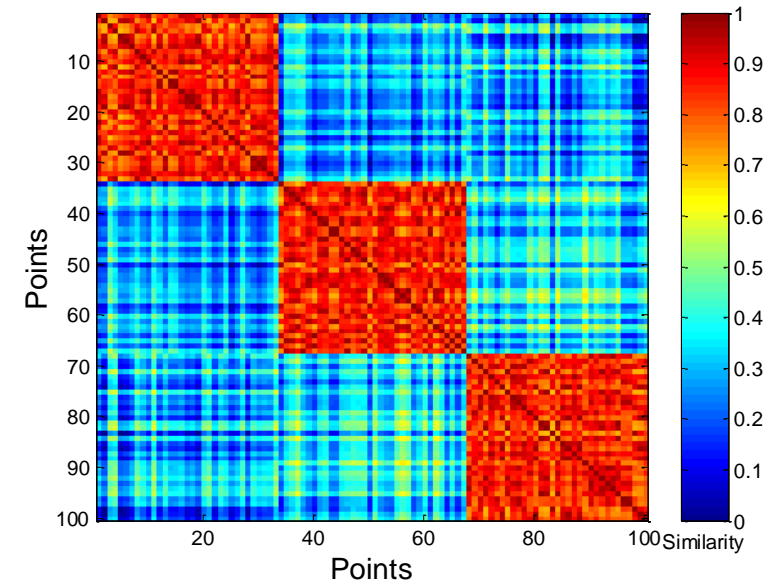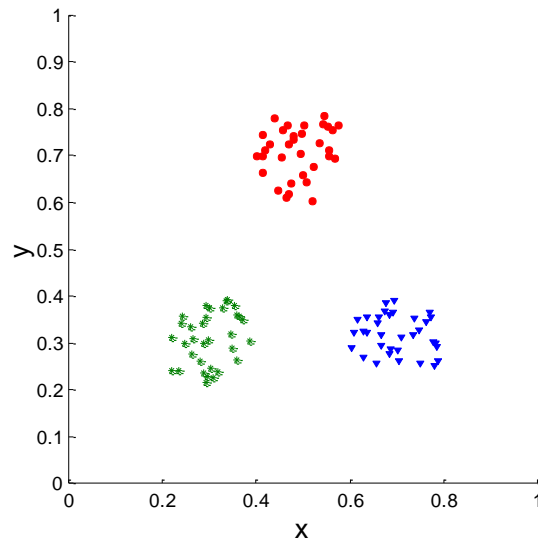


Corr = -0.9235



Corr = -0.5810

# Validity Measurement
## Using Similarity Matrix

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

- Clusters in random data are not so crisp
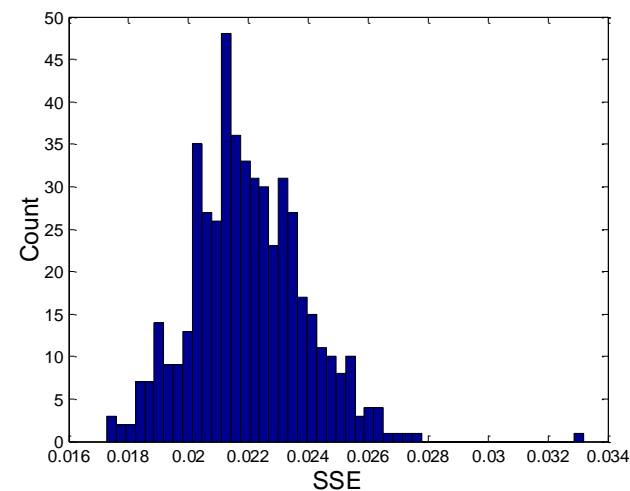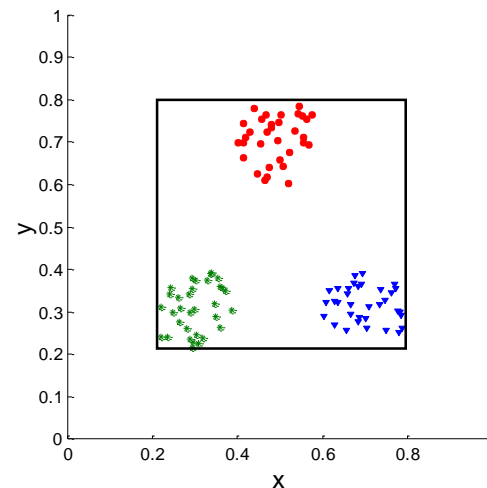
# Framework for
# Cluster Validity

- Need a framework to interpret any measure.
  - For example, if our measure of evaluation has the value, 10, is that good, fair, or poor?

- Statistics provide a framework for cluster validity
  - The more "atypical" a clustering result is, the more likely it represents valid structure in the data
  - Can compare the values of an index that result from random data or clusterings to those of a clustering result.
    - If the value of the index is unlikely, then the cluster results are valid
  - These approaches are more complicated and harder to understand.

- For comparing the results of two different sets of cluster analyses, a framework is less necessary.
  - However, there is the question of whether the difference between two index values is significant
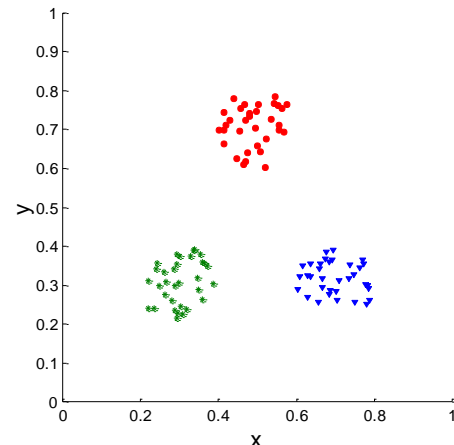
# Statistical Framework
## for SSE

**Example**

- Compare SSE of 0.005 against three clusters in random data

- Histogram shows SSE of three clusters in 500 sets of random data points of size 100 distributed over the range 0.018 – 0.028 for x and y values
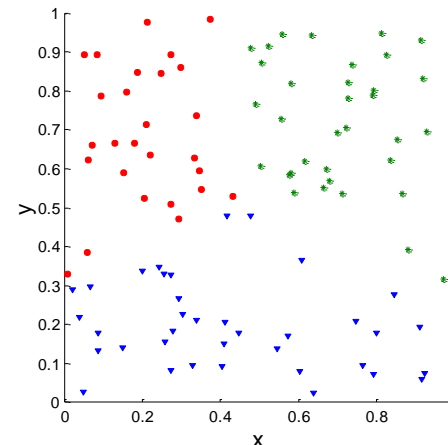
# Statistical Framework
## for Correlation

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



Corr = -0.9235



Corr = -0.5810