# Data Mining
# Cluster Analysis: Basic Concepts and Algorithms
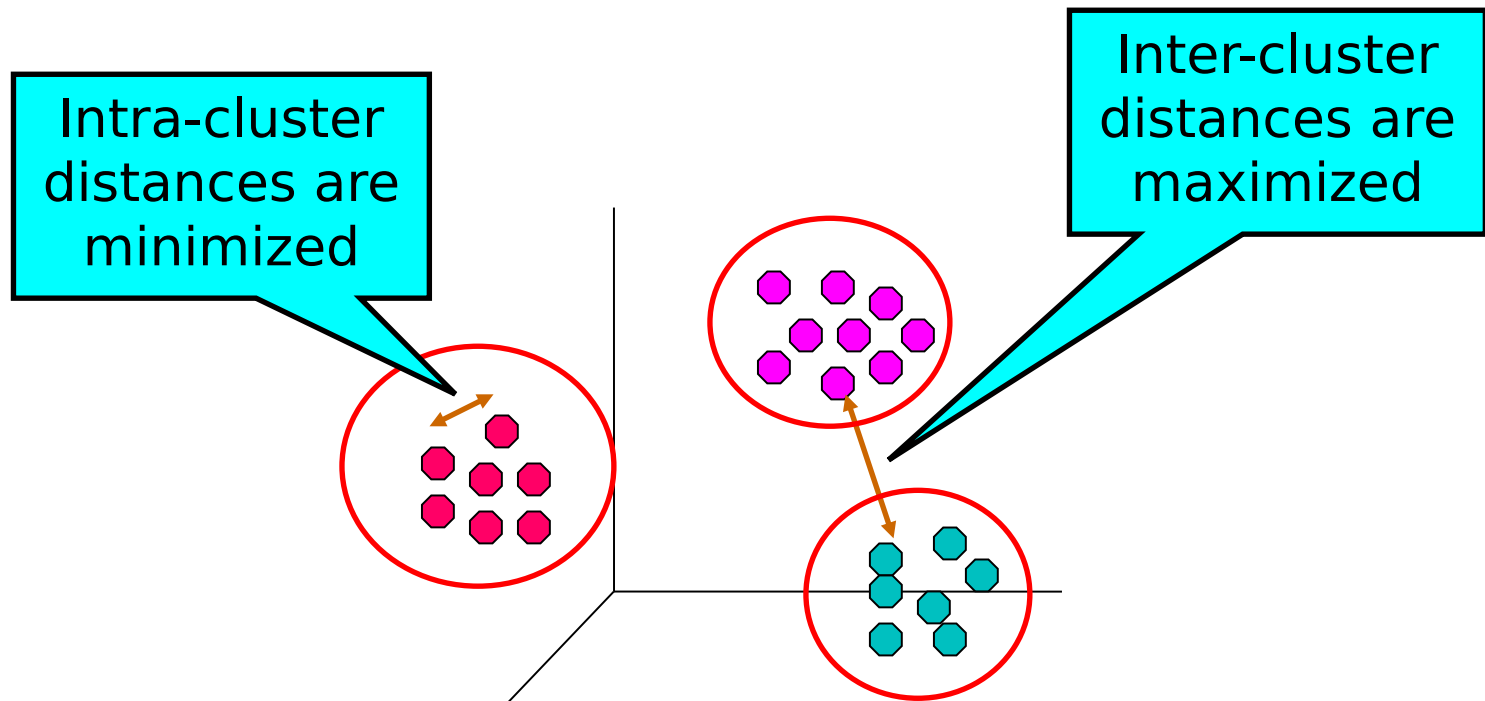
Lecture Notes for Chapter 7

Introduction to Data Mining
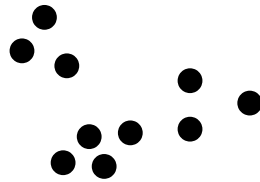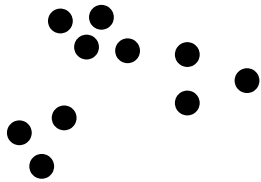
by

Tan, Steinbach, Kumar

# Clustering

- Unsupervised method
- Exploratory Data Analysis
- Useful for many applications like market segment analysis.
- What is clustering?

  – Organizing data into classes such that there is
    - High intra-class similarity
    - Low inter-class similarity

  } Finding the class labels and the number of classes directly from the data(in contrast to classification)

  } More informally, finding natural groupings among objects.

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**
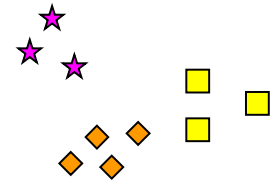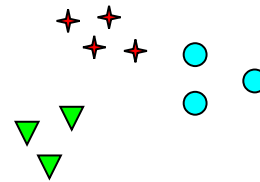
# What is Cluster Analysis?

- Given a set of objects, place them in groups such that the objects in a group are similar (or related) to one another and different from (or unrelated to) the objects in other groups

- The greater the similarity (or homogeneity) within a group and the greater the difference between groups, the better or more distinct the clustering.

Intra-cluster distances are minimized

Inter-cluster distances are maximized

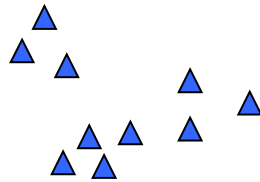**Introduction to Data Mining, 2nd Edition**
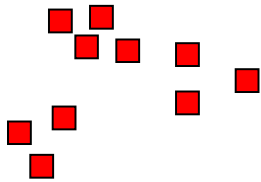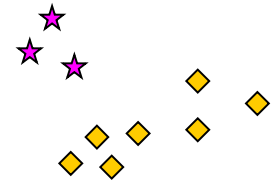**Tan, Steinbach, Karpatne, Kumar**

# Notion of a Cluster can be Ambiguous



How many clusters?

Six Clusters

Two Clusters

Four Clusters

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# Types of Clusterings

- A clustering is a set of clusters

  - An entire collection of clusters is commonly referred to as a clustering.

- Important distinction between hierarchical and partitional sets of clusters

  - Partitional Clustering

    - A division of data objects into **non-overlapping** subsets (clusters)

  - Hierarchical clustering

    - A set of **nested clusters** organized as a hierarchical tree

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# Partitional Clustering

**Original Points**

**A Partitional Clustering**

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# Hierarchical Clustering

**Traditional Hierarchical Clustering**

**Traditional Dendrogram**

**Non-traditional Hierarchical Clustering**

**Non-traditional Dendrogram**

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# Other Distinctions Between Sets of Clusters

- Exclusive versus non-exclusive
  - In non-exclusive clusterings, points may belong to multiple clusters.
    - Can belong to multiple classes or could be 'border' points
  - Fuzzy clustering  (one type of non-exclusive)
    - In fuzzy clustering, a point belongs to every cluster with some weight between 0 and 1
    - Weights must sum to 1
    - Probabilistic clustering has similar characteristics

- Partial versus complete
  - In some cases, we only want to cluster some of the data

# Types of Clusters

- Well-separated clusters

- Prototype-based clusters

- Contiguity-based clusters

- Density-based clusters

- Described by an Objective Function

# 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**

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# Types of Clusters: Prototype-Based

- Prototype-based (center-based)
  - A cluster is a set of objects such that an object in a cluster is closer (more similar) to the prototype or "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**

Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar

# Types of Clusters: Contiguity-Based

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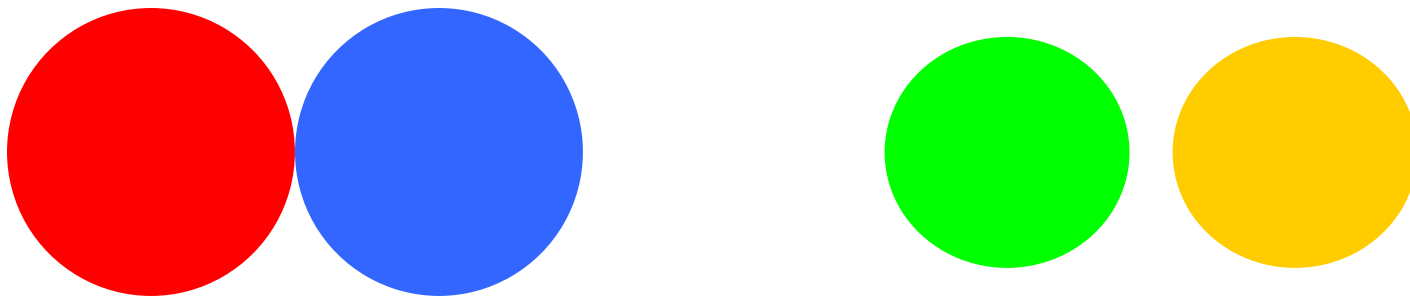  - A cluster is a set of points such that a point in a cluster is closer (or more similar) to one or more other points in the cluster than to any point not in the cluster.

**8 contiguous clusters**

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# 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**

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# 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
  - 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.

# Characteristics of the Input Data Are Important

- Type of proximity or density measure
  - Central to clustering
  - Depends on data and application

- Data characteristics that affect proximity and/or density are
  - Dimensionality
    - Sparseness
  - Attribute type
  - Special relationships in the data
    - For example, autocorrelation
  - Distribution of the data

- Noise and Outliers
  - Often interfere with the operation of the clustering algorithm
- Clusters of differing sizes, densities, and shapes

# Clustering Algorithms

- K-means and its variants

- Hierarchical clustering

- Density-based clustering

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# K-means Clustering

- Partitional clustering approach
- Number of clusters, K, must be specified
- Each cluster is associated with a centroid (center point)
- Each point is assigned to the cluster with the closest centroid
- 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 - Example

## Clustering Exercise

| X | Y |
|---|---|
| 2 | 4 |
| 2 | 6 |
| 5 | 6 |
| 4 | 7 |
| 8 | 3 |
| 6 | 6 |
| 5 | 2 |
| 5 | 7 |
| 6 | 3 |
| 4 | 4 |

**Introduction to Data Mining, 2nd Edition
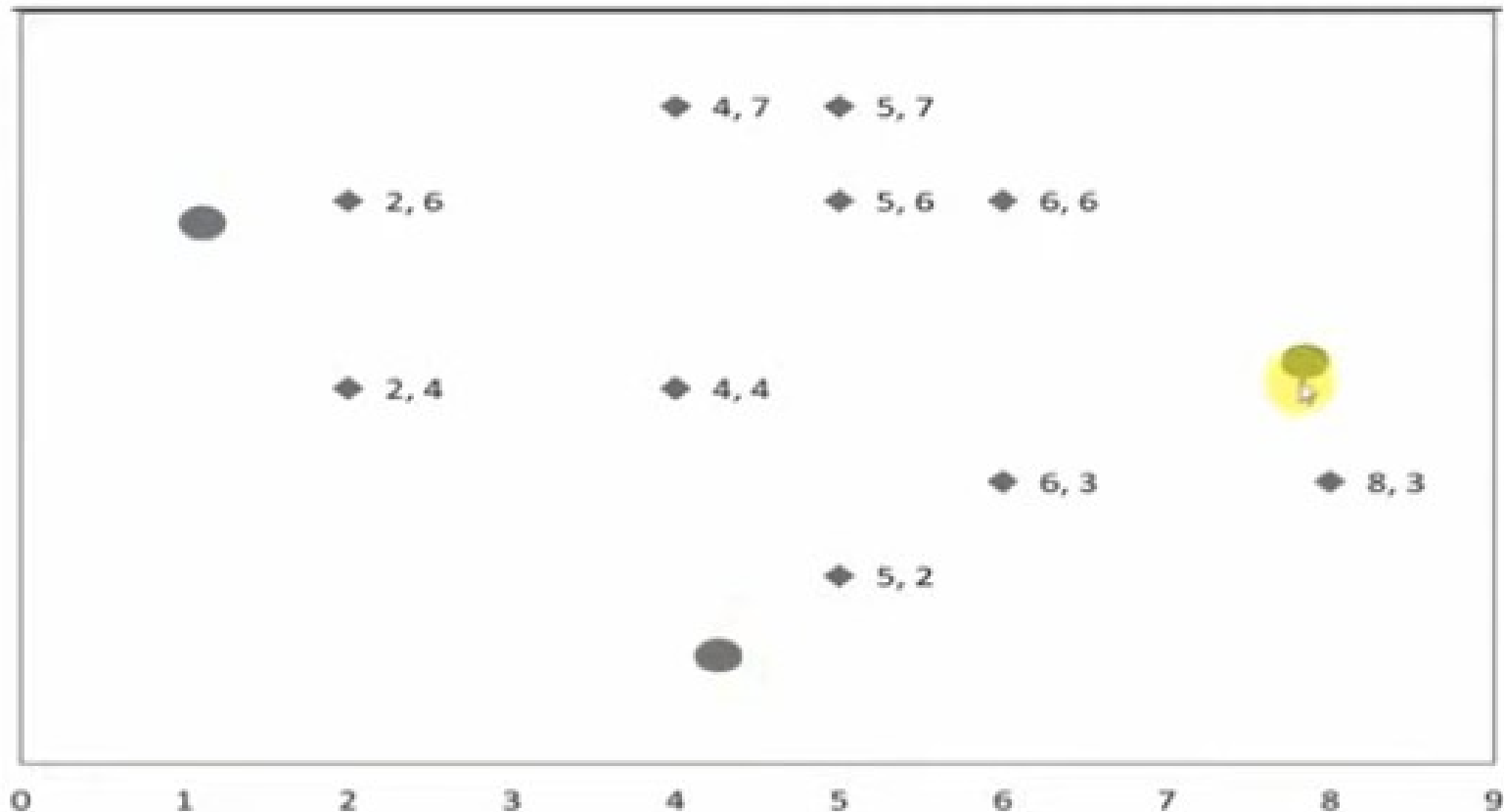Tan, Steinbach, Karpatne, Kumar**

# K-means Clustering - Example

- Initial case: All data points belongs to one cluster(with one centroid or seed point

# K-means Clustering - Example

- For making three clusters we need to identify three centroids
  - It can be data points or random points

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# K-means Clustering - Example

- For each data point distance is calculated from all the three centroids using some distance formula for ex. ED.
  - Assign the data points to cluster (minimum distance from centroid)

**Iteration - 1**

C1 - Seed Point1 – (1, 5)
C2 - Seed Point2 – (4, 1)
C3 - Seed Point3 – ( 8, 4)

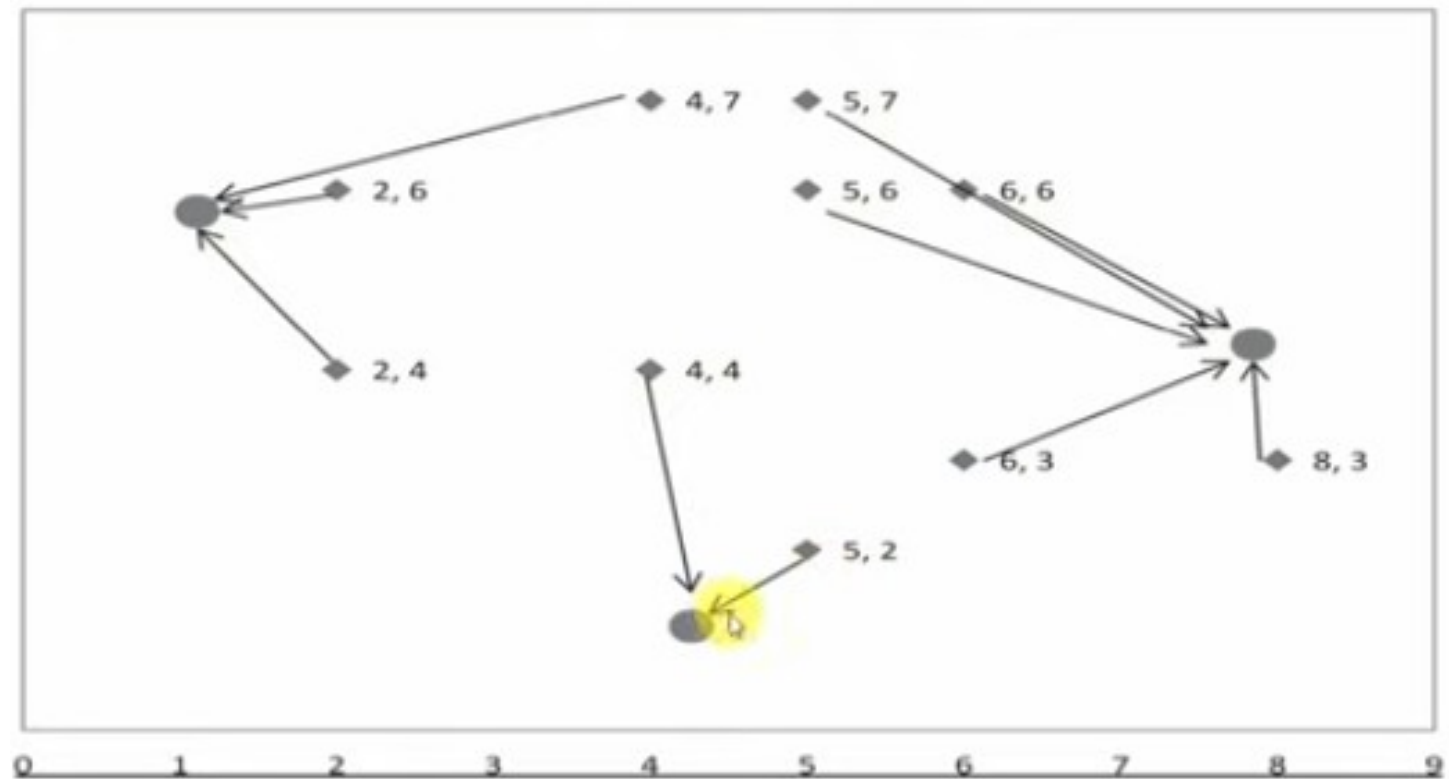$$D = \sqrt{((x_2 - x_1)^2 + (y_2 - y_1)^2)}$$

C1 – Centroid – (2.66, 5.66)
C2 – Centroid – (4.5, 3)
C3 – Centroid – ( 6, 5)

| | | | Distance to | | Cluster |
| X | Y | (1, 5) | (4, 1) | (8, 4) | Number |
|---|---|---|---|---|---|
| 2 | 4 | 1.41 | 3.61 | 6.00 | C1 |
| 2 | 6 | 1.41 | 5.39 | 6.32 | C1 |
| 5 | 6 | 4.12 | 5.10 | 3.61 | C3 |
| 4 | 7 | 3.61 | 6.00 | 5.00 | C1 |
| 8 | 3 | 7.28 | 4.47 | 1.00 | C3 |
| 6 | 6 | 5.10 | 5.39 | 2.83 | C3 |
| 5 | 2 | 5.00 | 1.41 | 3.61 | C2 |
| 5 | 7 | 4.47 | 6.08 | 4.24 | C3 |
| 6 | 3 | 5.39 | 2.83 | 2.24 | C3 |
| 4 | 4 | 3.16 | 3.00 | 4.00 | C2 |

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# K-means Clustering - Example

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# K-means Clustering - Example

- Previously all data points were in <u>one cluster</u>

- Now data points are in <u>three cluseters(C1,C2,C3)</u>

- So there is <u>movement of data points</u> from one cluster to three different clusters.

- Now we need to <u>again calculate the distance of all data points with new centroids</u>

- New Centroids are calculated as:

    - **(Average of X coordinates, Average of Y coordinates)**

    - For example new Centroid for C1 is:

        - **(2+2+4)/3 , (4+6+7)/3 =(2.66, 5.66)**

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# K-means Clustering - Example

**Iteration - 2**
C1 – Centroid – (2.66, 5.66)
C2 – Centroid – (4.5, 3)
C3 – Centroid – ( 6, 5)


C1 – Centroid – (2.66, 5.66)
C2 – Centroid – (5, 3)
C3 – Centroid – ( 6, 5.5)

| | | Distance to | | | Cluster |
| X | Y | (2.66, 5.66) | (4.5, 3) | (6, 5) | Number |
|---|---|---|---|---|---|
| 2 | 4 | 1.79 | 2.69 | 4.12 | C1 |
| 2 | 6 | 0.74 | 3.91 | 4.12 | C1 |
| 5 | 6 | 2.36 | 3.04 | 1.41 | C3 |
| 4 | 7 | 1.90 | 4.03 | 2.83 | C1 |
| 8 | 3 | 5.97 | 3.5 | 2.83 | C3 |
| 6 | 6 | 3.36 | 3.35 | 1 | C3 |
| 5 | 2 | 4.34 | 1.12 | 3.16 | C2 |
| 5 | 7 | 2.70 | 4.03 | 2.24 | C3 |
| 6 | 3 | 4.27 | 1.5 | 2 | C2 |
| 4 | 4 | 2.13 | 1.12 | 2.24 | C2 |

# K-means Clustering - Example

**Iteration - 3**

C1 – Centroid – (2.66, 5.66)
C2 – Centroid – (5, 3)
C3 – Centroid – ( 6, 5.5)

C1 – Centroid – (2.66, 5.66)
C2 – Centroid – (5.75, 3)
C3 – Centroid – ( 5.33, 6.33)

| | | | Distance to | | | Cluster |
|---|---|---|---|---|---|---|
| X | Y | (2.66, 5.66) | (5, 3) | (6, 5.5) | Number |
| 2 | 4 | 1.79 | 3.16 | 4.27 | C1 |
| 2 | 6 | 0.74 | 4.24 | 4.03 | C1 |
| 5 | 6 | 2.36 | 3.00 | 1.12 | C3 |
| 4 | 7 | 1.90 | 4.12 | 2.50 | C1 |
| 8 | 3 | 5.97 | 3.00 | 3.20 | C2 |
| 6 | 6 | 3.36 | 3.16 | 0.50 | C3 |
| 5 | 2 | 4.34 | 1.00 | 3.64 | C2 |
| 5 | 7 | 2.70 | 4.00 | 1.80 | C3 |
| 6 | 3 | 4.27 | 1.00 | 2.50 | C2 |
| 4 | 4 | 2.13 | 1.41 | 2.50 | C2 |

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# K-means Clustering - Example

**Iteration - 4**

C1 – Centroid – (2.66, 5.66)
C2 – Centroid – (5.75, 3)
C3 – Centroid – ( 5.33, 6.33)

C1 – Centroid – (2, 5)
C2 – Centroid – (5.75, 3)
C3 – Centroid – ( 5, 6.5)

|   |   | Distance to | | | Cluster |
| X | Y | (2.66, 5.66) | (5.75, 3) | (5.33, 6.33) | Number |
|---|---|---|---|---|---|
| 2 | 4 | 1.79 | 3.88 | 4.06 | C1 |
| 2 | 6 | 0.74 | 4.80 | 3.35 | C1 |
| 5 | 6 | 2.36 | 3.09 | 0.47 | C3 |
| 4 | 7 | 1.90 | 4.37 | 1.49 | C3 |
| 8 | 3 | 5.97 | 2.25 | 4.27 | C2 |
| 6 | 6 | 3.36 | 3.01 | 0.75 | C3 |
| 5 | 2 | 4.34 | 1.25 | 4.34 | C2 |
| 5 | 7 | 2.70 | 4.07 | 0.75 | C3 |
| 6 | 3 | 4.27 | 0.25 | 3.40 | C2 |
| 4 | 4 | 2.13 | 2.02 | 2.68 | C2 |

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# K-means Clustering - Example

**Iteration - 5**

C1 – Centroid – (2, 5)
C2 – Centroid – (5.75, 3)
C3 – Centroid – ( 5, 6.5)

No movement of data Points
Hence these are the final
positions

| | | | Distance to | | Cluster |
|---|---|---|---|---|---|
| X | Y | (2, 5) | (5.75, 3) | (5, 6.5) | Number |
| 2 | 4 | 1.00 | 3.88 | 3.91 | C1 |
| 2 | 6 | 1.00 | 4.80 | 3.04 | C1 |
| 5 | 6 | 3.16 | 3.09 | 0.50 | C3 |
| 4 | 7 | 2.83 | 4.37 | 1.12 | C3 |
| 8 | 3 | 6.32 | 2.25 | 4.61 | C2 |
| 6 | 6 | 4.12 | 3.01 | 1.12 | C3 |
| 5 | 2 | 4.24 | 1.25 | 4.50 | C2 |
| 5 | 7 | 3.61 | 4.07 | 0.50 | C3 |
| 6 | 3 | 4.47 | 0.25 | 3.64 | C2 |
| 4 | 4 | 2.24 | 2.02 | 2.69 | C2 |

# K-means Clustering - Example

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# K-means Clustering – Details

- Simple iterative algorithm.
  - Choose initial centroids;
  - repeat {assign each point to a nearest centroid; re-compute cluster centroids}
  - until centroids stop changing.

- Initial centroids are often chosen randomly.
  - Clusters produced can vary from one run to another

- The centroid is (typically) the mean of the points in the cluster, but other definitions are possible (see Table 7.2).

- K-means will converge for common proximity measures with appropriately defined centroid (see Table 7.2)

- 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

# K-means Objective Function

- A common objective function (used with Euclidean distance measure) is Sum of Squared Error (SSE)
  - For each point, the error is the distance to the nearest cluster center
  - 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 $C_i$ and $m_i$ is the centroid (mean) for cluster $C_i$
  - SSE improves in each iteration of K-means until it reaches a local or global minima.

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Two different K-means Clusterings



Original Points

Optimal Clustering

Sub-optimal Clustering

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

Iteration 5

# Importance of Choosing Initial Centroids ...

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# Importance of Choosing Intial Centroids

- Depending on the choice of initial centroids, B and C may get merged or remain separate

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# 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 = $10!/10^{10}$ = 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

# 10 Clusters Example
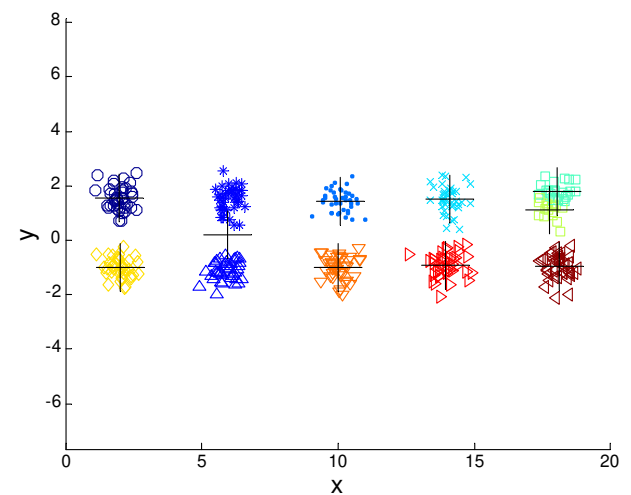


Iteration 4

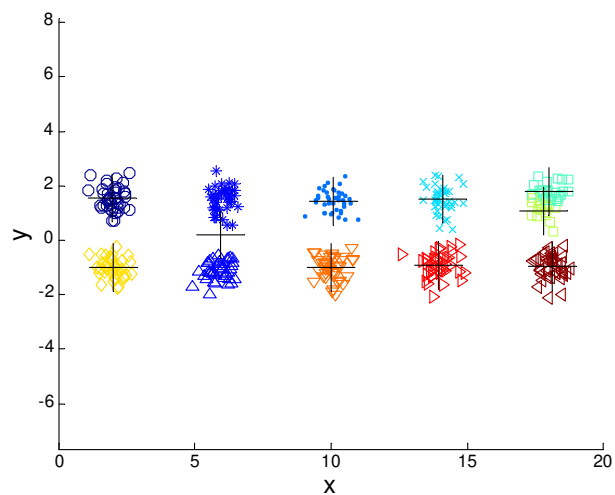**Starting with two initial centroids in one cluster of each pair of clusters**
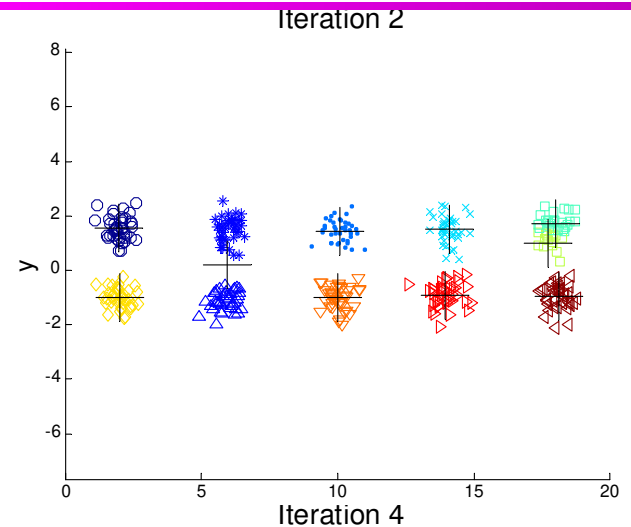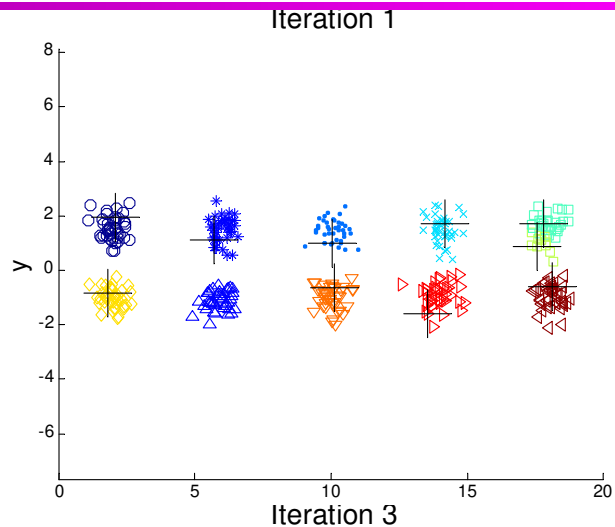
# 10 Clusters Example



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

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# 10 Clusters Example



Iteration 4

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

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# 10 Clusters Example



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

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# Solutions to Initial Centroids Problem

- Multiple runs
  - Helps, but probability is not on your side
- Use some strategy to select the k initial centroids and then select among these initial centroids
  - Select most widely separated
    - K-means++ is a robust way of doing this selection
  - Use hierarchical clustering to determine initial centroids
- Bisecting K-means
  - Not as susceptible to initialization issues

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# K-means++

- This approach can be slower than random initialization, but very consistently produces better results in terms of SSE

    – The k-means++ algorithm guarantees an approximation ratio O(log k) in expectation, where k is the number of centers

- To select a set of initial centroids, *C*, perform the following

1. Select an initial point at random to be the first centroid

2. For k – 1 steps

3. For each of the N points, $x_i$, $1 \leq i \leq N$, find the minimum squared distance to the currently selected centroids, $C_1, \ldots, C_{j,} 1 \leq j < k$, i.e.,

4. Randomly select a new centroid by choosing a point with probability proportional to is

5. End For

# Bisecting K-means

- ## Bisecting K-means algorithm
  - Variant of K-means that can produce a partitional or a hierarchical clustering

**Algorithm 3** Bisecting K-means Algorithm.

1: Initialize the list of clusters to contain the cluster containing all points.
2: **repeat**
3:   Select a cluster from the list of clusters
4:   **for** $i = 1$ to $number\_of\_iterations$ **do**
5:     Bisect the selected cluster using basic K-means
6:   **end for**
7:   Add the two clusters from the bisection with the lowest SSE to the list of clusters.
8: **until** Until the list of clusters contains $K$ clusters

**CLUTO: http://glaros.dtc.umn.edu/gkhome/cluto/cluto/overview**

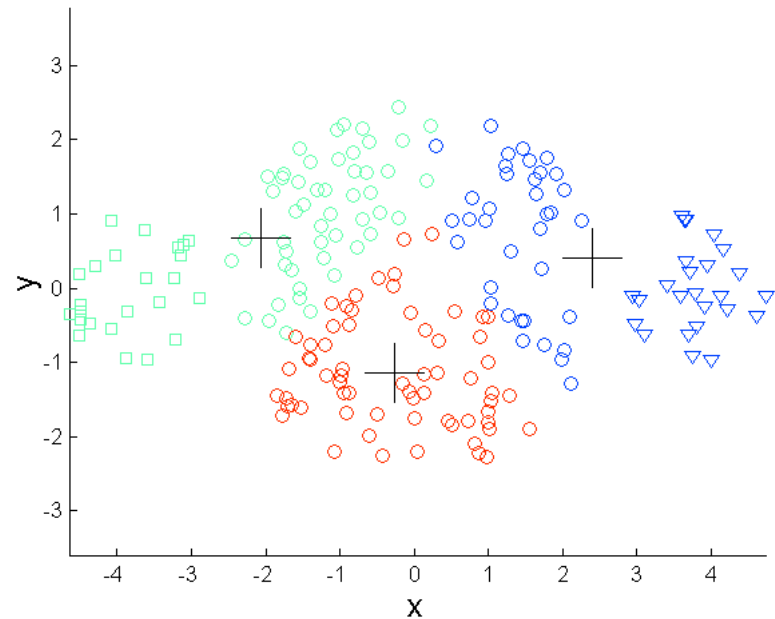# Bisecting K-means Example

# 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.
  - One possible solution is to remove outliers before clustering
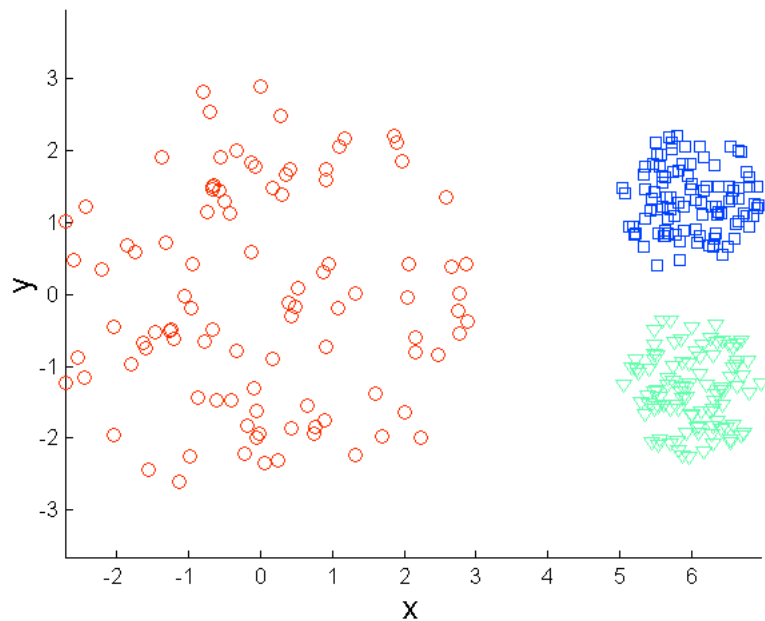
**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Limitations of K-means: Differing Sizes



**Original Points**

**K-means (3 Clusters)**

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Limitations of K-means: Differing Density



**Original Points**

**K-means (3 Clusters)**

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# Limitations of K-means: Non-globular Shapes



**Original Points**

**K-means (2 Clusters)**

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**
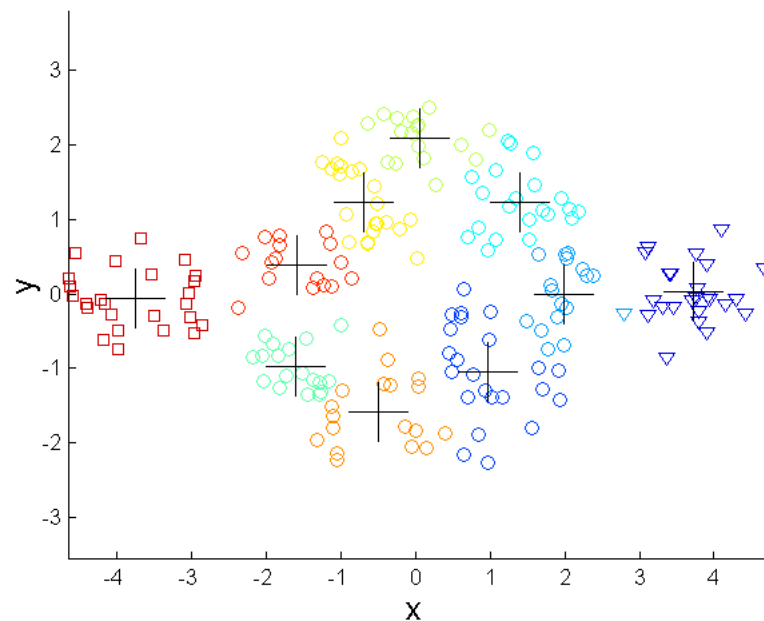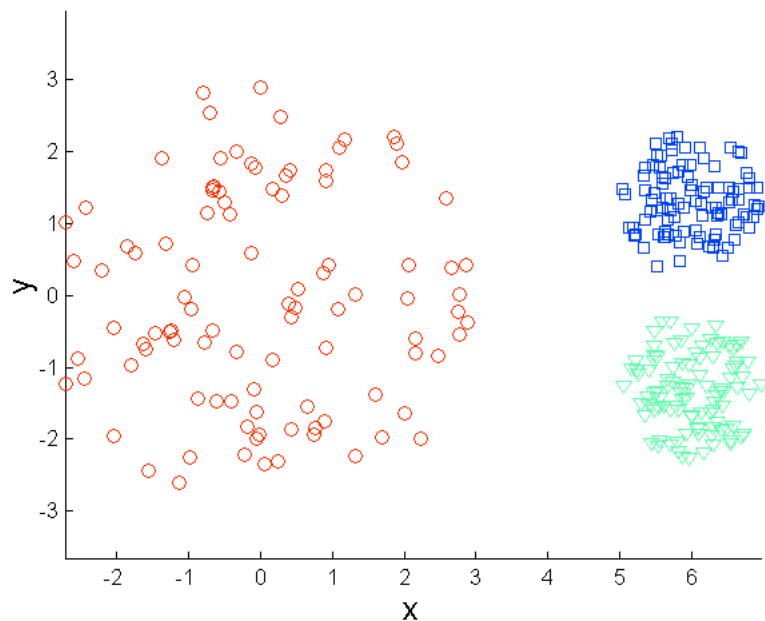
# Overcoming K-means Limitations
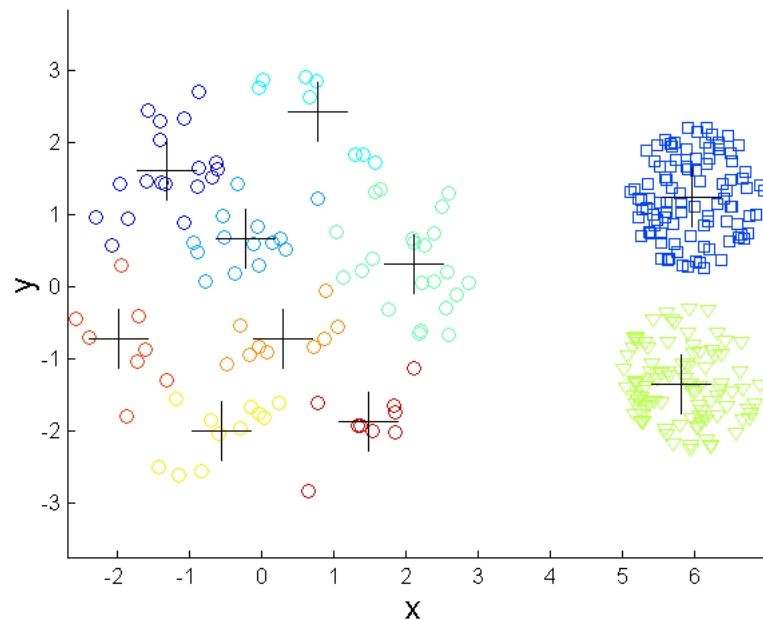


**Original Points**



**K-means Clusters**

One solution is to find a large number of clusters such that each of them represents a part of a natural cluster. But these small clusters need to be put together in a post-processing step.

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**
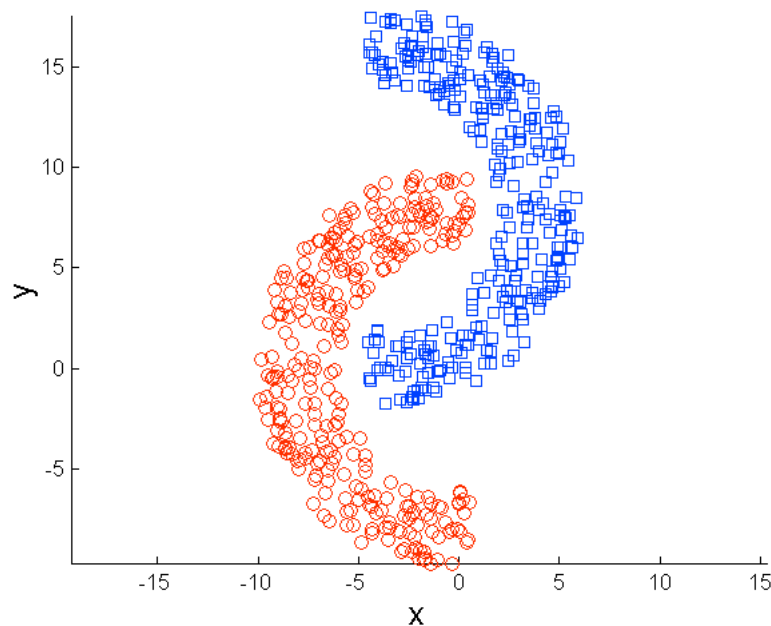
# Overcoming K-means Limitations
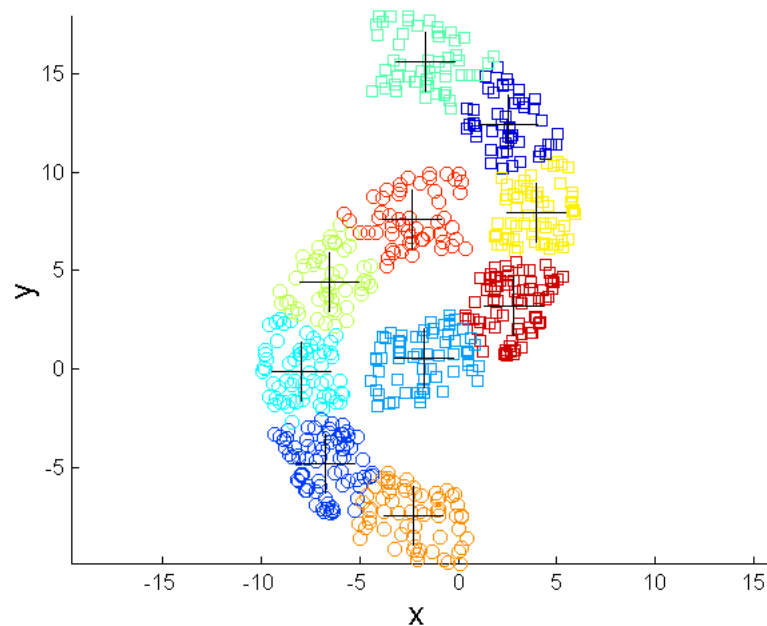


**Original Points**

**K-means Clusters**

One solution is to find a large number of clusters such that each of them represents a part of a natural cluster. But these small clusters need to be put together in a post-processing step.

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Overcoming K-means Limitations
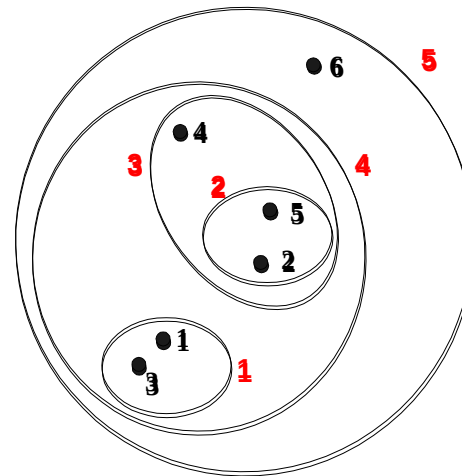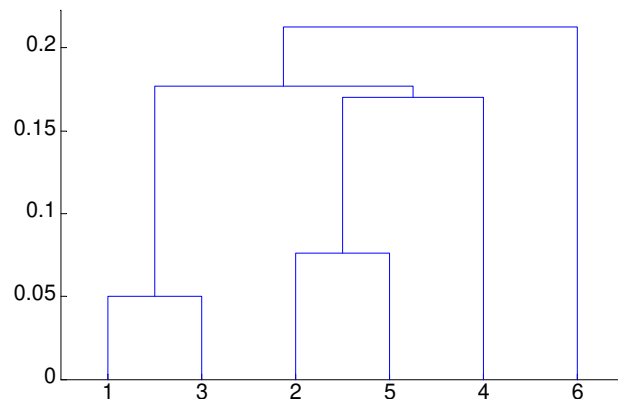


**Original Points**

**K-means Clusters**

One solution is to find a large number of clusters such that each of them represents a part of a natural cluster. But these small clusters need to be put together in a post-processing step.

# 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, …)
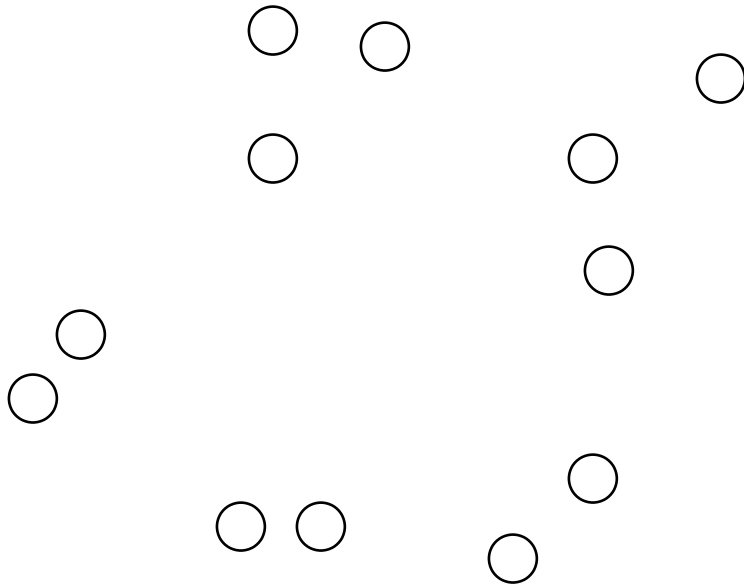
# Hierarchical Clustering

- 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 an individual 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

- ### Key Idea: Successively merge closest clusters

- Basic algorithm
  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

# Steps 1 and 2

- 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

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Intermediate Situation

- After some merging steps, we have some clusters



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

**Proximity Matrix**

p1  p2   p3  p4   …   p9   p10  p11  p12
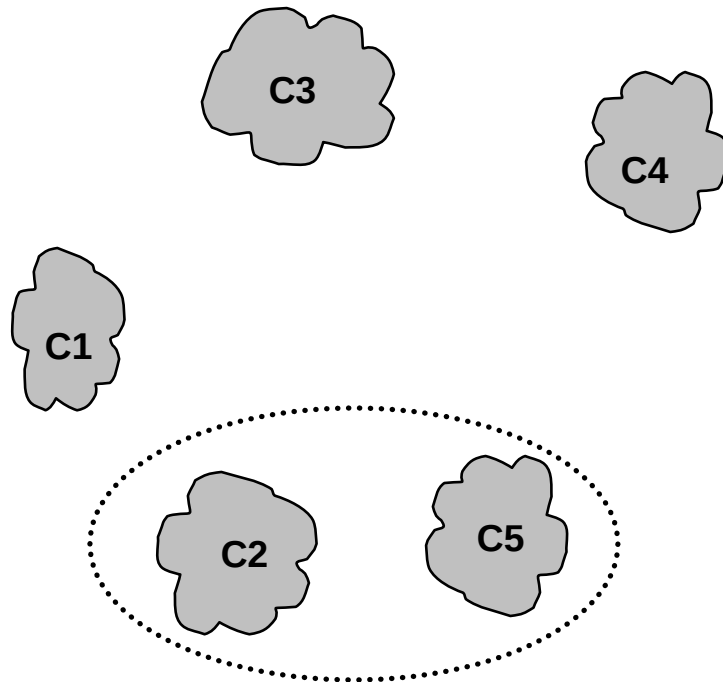
**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Step 4

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



**Proximity Matrix**

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Step 5

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



**Proximity Matrix**

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

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# How to Define Inter-Cluster Distance

**Similarity?**

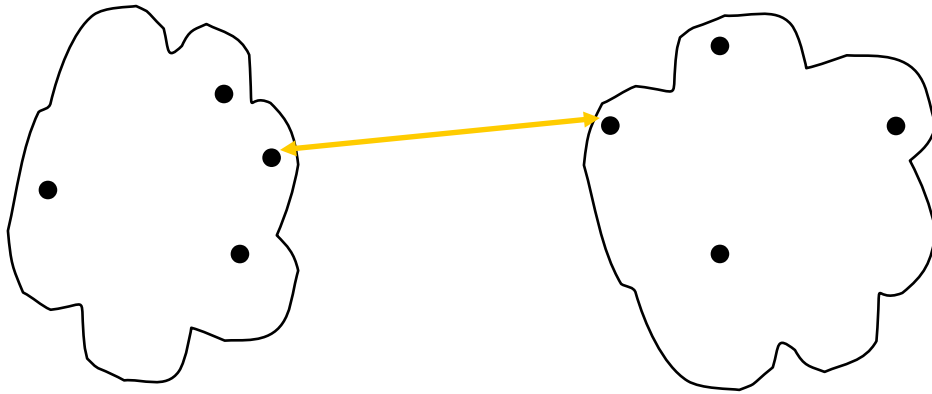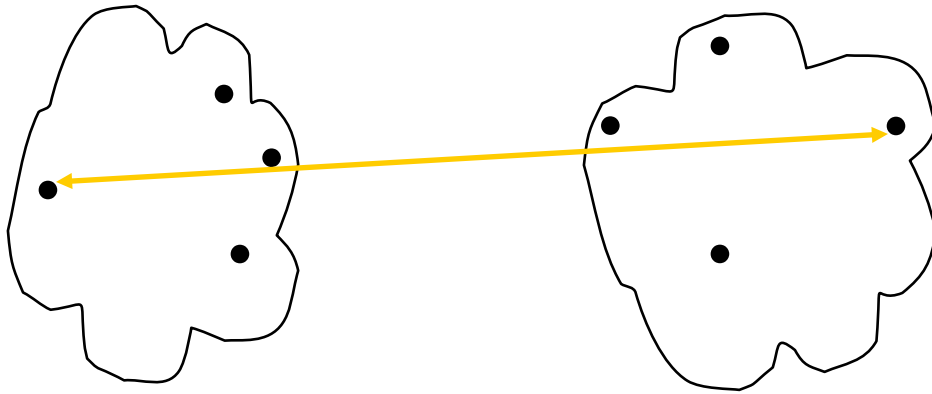|  | p1 | p2 | p3 | p4 | p5 | . . . |
|---|---|---|---|---|---|---|
| **p1** | | | | | | |
| **p2** | | | | | | |
| **p3** | | | | | | |
| **p4** | | | | | | |
| **p5** | | | | | | |
| **.** | | | | | | |
| **.** | | | | | | |
| **.** | | | | | | |

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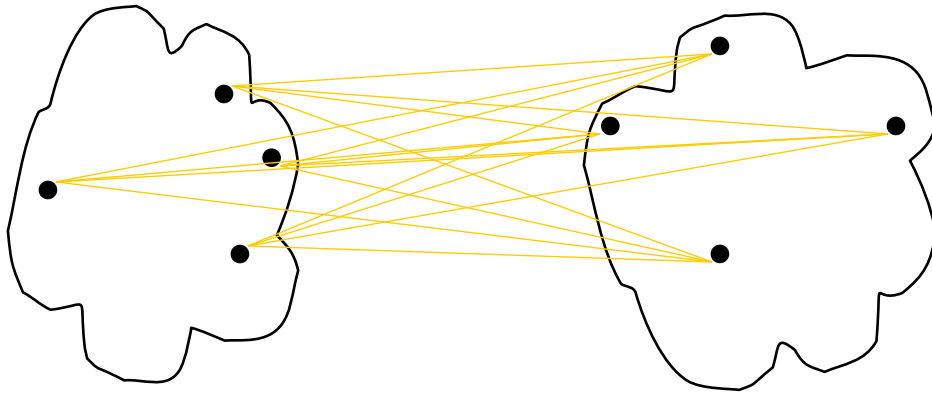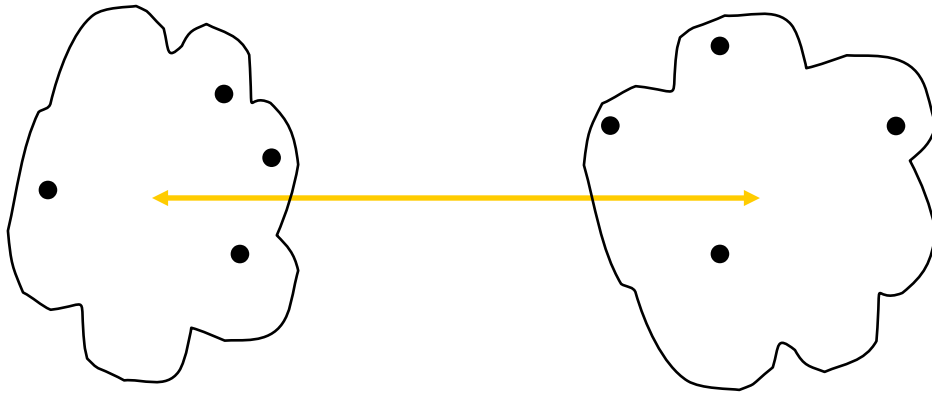
- <span style="color:red">MIN</span>
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
    - Ward's Method uses squared error

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# How to Define Inter-Cluster Similarity

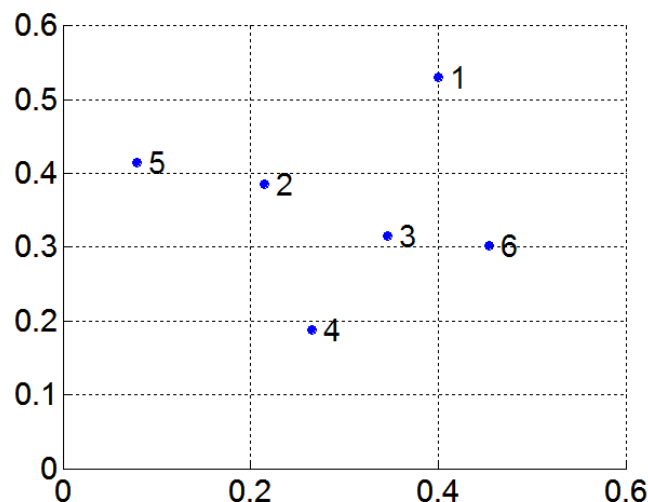|    | p1 | p2 | p3 | p4 | p5 | . . . |
|----|----|----|----|----|----|-------|
| p1 |    |    |    |    |    |       |
| p2 |    |    |    |    |    |       |
| p3 |    |    |    |    |    |       |
| p4 |    |    |    |    |    |       |
| p5 |    |    |    |    |    |       |
| .  |    |    |    |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

**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

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# How to Define Inter-Cluster Similarity



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

**Proximity Matrix**

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

# MIN or Single Link

- Proximity of two clusters is based on the two closest points in the different clusters
  - Determined by one pair of points, i.e., by one link in the proximity graph

- 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**

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Strength of MIN



Original Points

Six Clusters

**Can handle non-elliptical shapes**

# Limitations of MIN

**Two Clusters**

**Original Points**
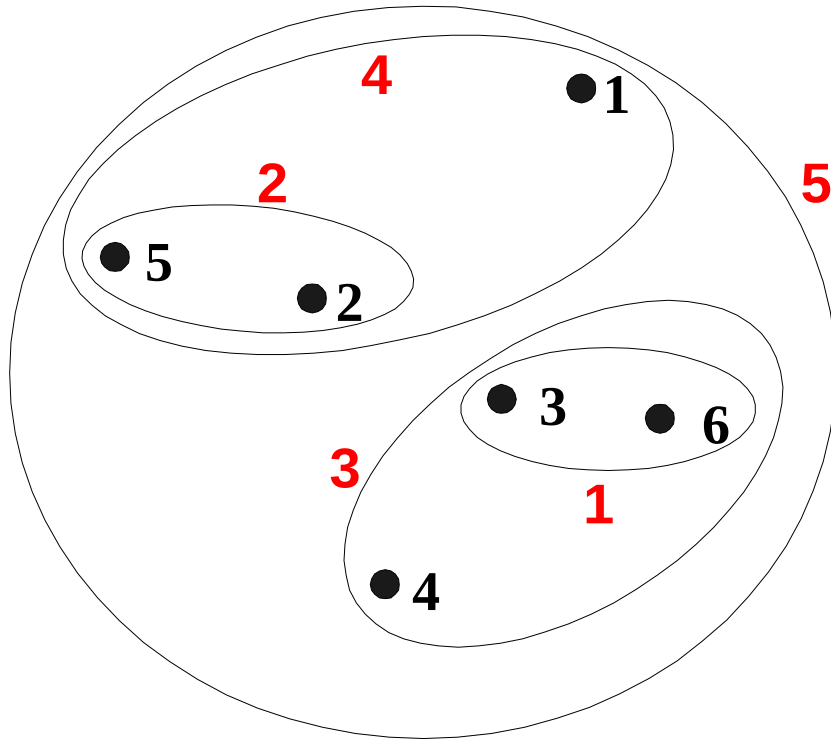
**Sensitive to noise**

**Three Clusters**

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# MAX or Complete Linkage

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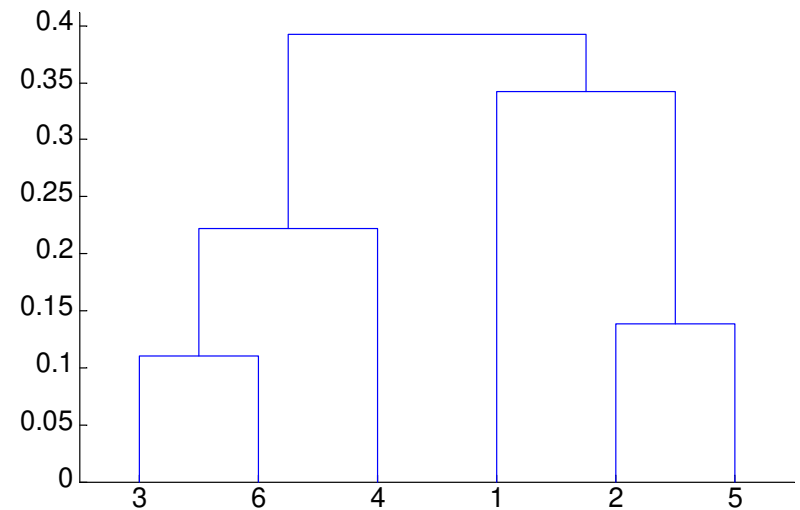  - Determined by all pairs of points in the two clusters

**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**

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Strength of MAX



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

**Less susceptible to noise**

Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar

# Limitations of MAX



**Original Points**

**Two Clusters**

**Tends to break large clusters**

**Biased towards globular clusters**

Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar

# Group Average

- Proximity of two clusters is the average of pairwise proximity between points in the two clusters.
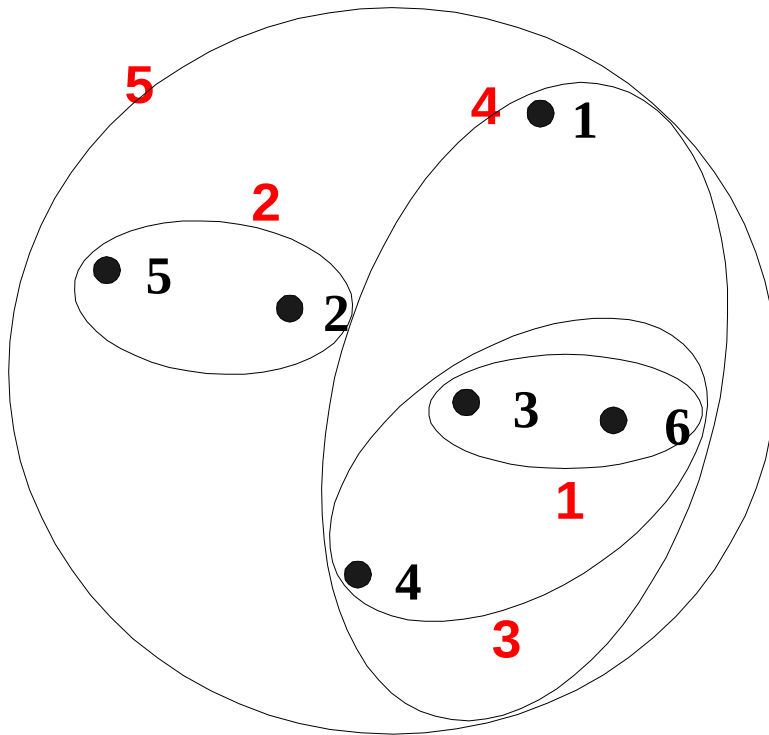
$$\text{proximity}(\text{Cluster}_i, \text{Cluster}_j) = \frac{\sum\limits_{\substack{p_i \in \text{Cluster}_i \\ p_j \in \text{Cluster}_j}} \text{proximity}(p_i, p_j)}{|\text{Cluster}_i| \times |\text{Cluster}_j|}$$
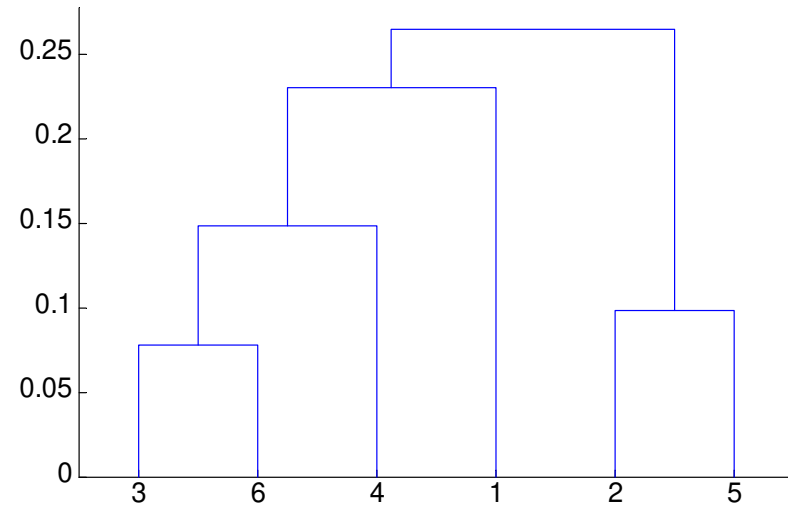
**Distance Matrix:**

|    | 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**

**Introduction to Data Mining, 2nd Edition
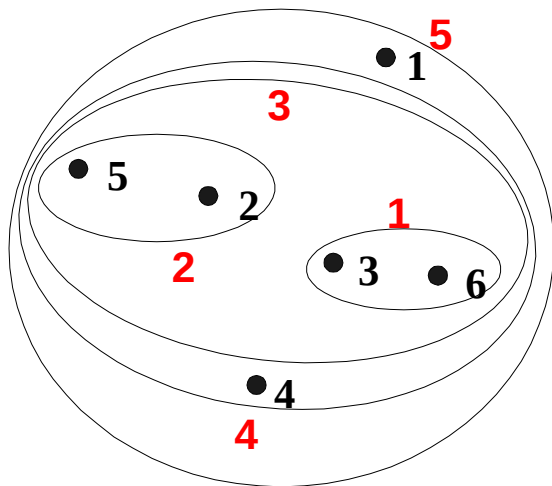Tan, Steinbach, Karpatne, Kumar**

# Hierarchical Clustering: Group Average

- Compromise between Single and Complete Link

- Strengths
  - Less susceptible to noise

- Limitations
  - Biased towards globular clusters

**Introduction to Data Mining, 2nd Edition
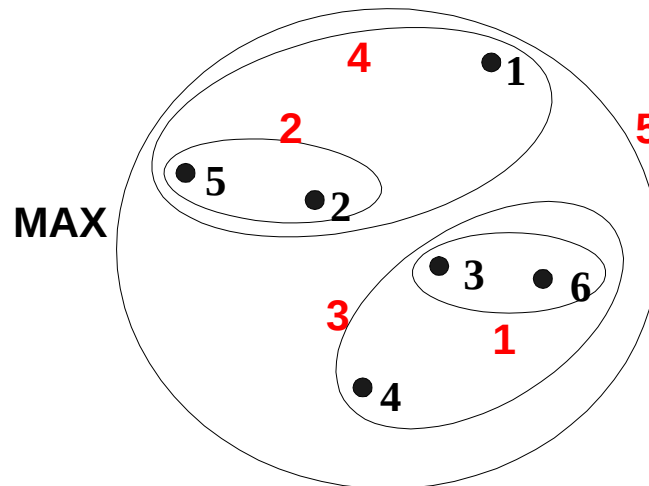Tan, Steinbach, Karpatne, Kumar**

# 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

- Biased towards globular clusters

- Hierarchical analogue of K-means
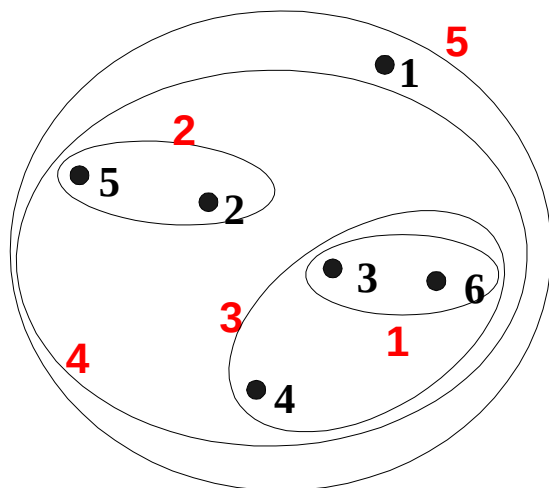  - Can be used to initialize K-means

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# Hierarchical Clustering: Comparison



**MIN**

**MAX**

**Group Average**

**Ward's Method**

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**
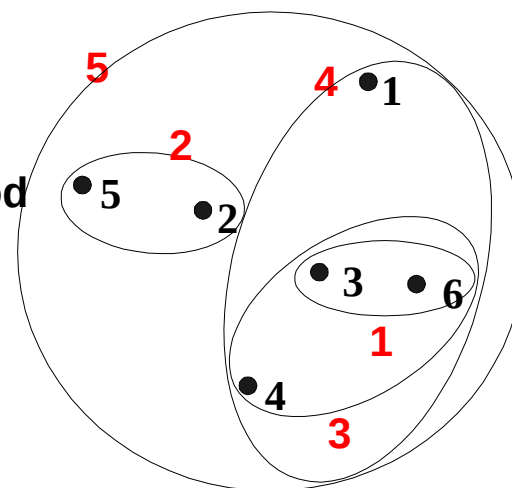
# Hierarchical Clustering: Time and Space requirements

- $O(N^2)$ 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 with some cleverness

# 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 clusters of different sizes and non-globular shapes
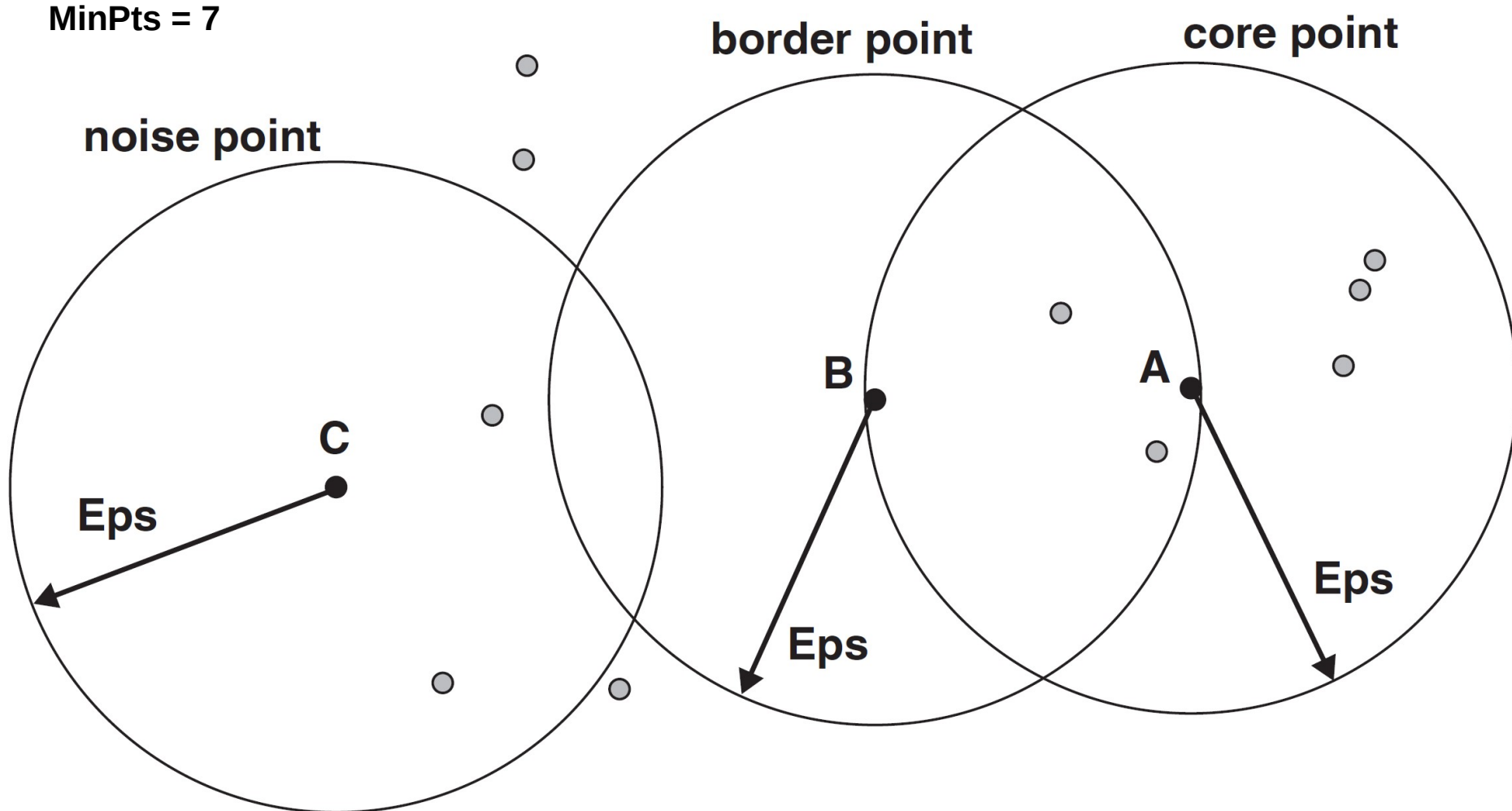  - Breaking large clusters

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# Density Based Clustering

- Clusters are regions of high density that are separated from one another by regions on low density.

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# DBSCAN

- DBSCAN is a density-based algorithm.
  - Density = number of points within a specified radius (Eps)

  - A point is a <span style="color:red">core point</span> if it has at least a specified number of points (MinPts) within Eps
    - These are points that are at the interior of a cluster
    - Counts the point itself

  - A <span style="color:red">border point</span> is not a core point, but is in the neighborhood of a core point

  - A <span style="color:red">noise point</span> is any point that is not a core point or a border point

# DBSCAN: Core, Border, and Noise Points



MinPts = 7

noise point
border point
core point

B
A
C
Eps
Eps
Eps

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# DBSCAN: Core, Border and Noise Points



**Original Points**

**Point types: core,
border and noise**

**Eps = 10, MinPts = 4**

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# DBSCAN Algorithm

- Form clusters using core points, and assign border points to one of its neighboring clusters

1: Label all points as core, border, or noise points.

2: Eliminate noise points.

3: Put an edge between all core points within a distance *Eps* of each other.

4: Make each group of connected core points into a separate cluster.

5: Assign each border point to one of the clusters of its associated core points

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# When DBSCAN Works Well



**Original Points**



**Clusters (**dark blue points indicate noise**)**

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

**Resistant to noise**

Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar

# When DBSCAN Does NOT Work Well



**Original Points**

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# When DBSCAN Does NOT Work Well



(MinPts=4, Eps=9.92).

**Original Points**



**Varying densities**

**High-dimensional data**

(MinPts=4, Eps=9.75)

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# DBSCAN: Determining EPS and MinPts

- Idea is that for points in a cluster, their $k^{th}$ nearest neighbors are at close distance

- Noise points have the $k^{th}$ nearest neighbor at farther distance
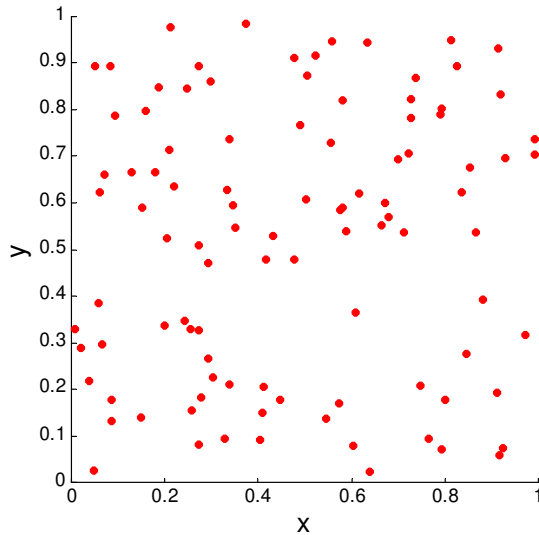
- So, plot sorted distance of every point to its $k^{th}$ nearest neighbor

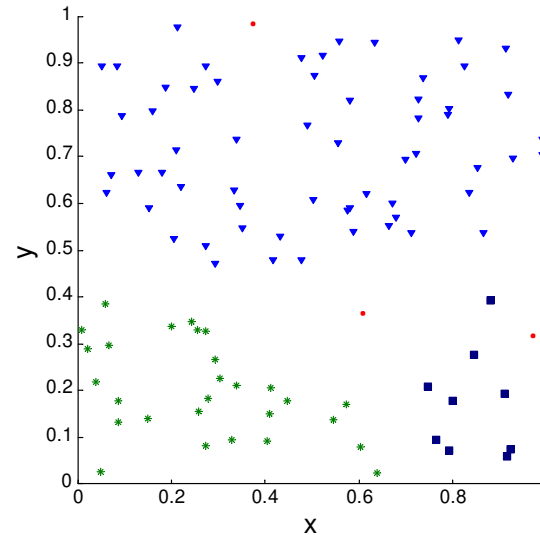**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# 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"!
  - In practice the clusters we find are defined by the clustering algorithm

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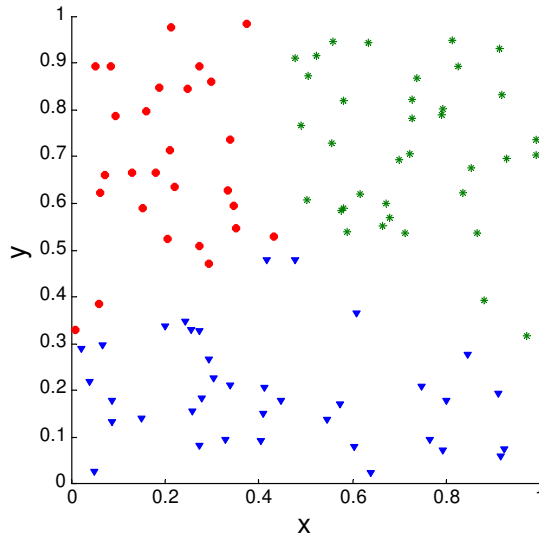  - To compare two clusters

**Introduction to Data Mining, 2nd Edition
Tan, Steinbach, Karpatne, Kumar**

# Clusters found in Random Data



**Random Points**

**DBSCAN**

**K-means**

**Complete Link**
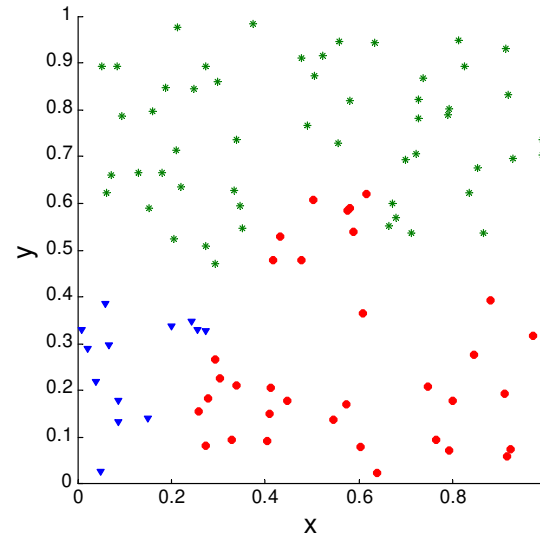
**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Measures of Cluster Validity

- Numerical measures that are applied to judge various aspects of cluster validity, are classified into the following two types.

    - Supervised: Used to measure the extent to which cluster labels match externally supplied class labels.
        - Entropy
        - Often called *external indices* because they use information external to the data

    - Unsupervised:  Used to measure the goodness of a clustering structure *without* respect to external information.
        - Sum of Squared Error (SSE)
        - Often called *internal indices* because they only use information in the data

- You can use supervised or unsupervised measures to compare clusters or clusterings

**Introduction to Data Mining, 2nd Edition
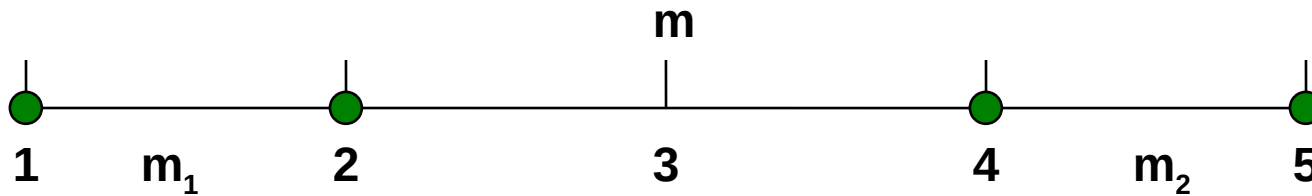Tan, Steinbach, Karpatne, Kumar**

# Unsupervised 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)

$$SSE = \sum_{i} \sum_{x \in C_i} \left( x - m_i \right)^2$$

  - Separation is measured by the between cluster sum of squares

$$SSB = \sum_{i} \left| C_i \right| \left( m - m_i \right)^2$$

  Where is the size of cluster $i$

# Unsupervised Measures: Cohesion and Separation

- Example: SSE
  - SSB + SSE = constant

**m**



| 1 | m$_1$ | 2 | 3 | 4 | m$_2$ | 5 |

**K=1 cluster:**

$$SSE = (1-3)^2 + (2-3)^2 + (4-3)^2 + (5-3)^2 = 10$$
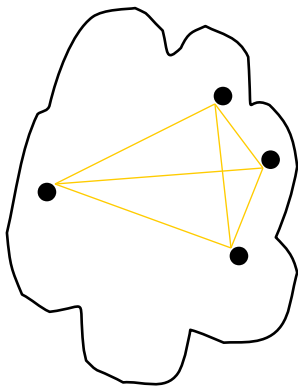
$$SSB = 4 \times (3-3)^2 = 0$$

$$Total = 10 + 0 = 10$$

**K=2 clusters:**

$$SSE = (1-1.5)^2 + (2-1.5)^2 + (4-4.5)^2 + (5-4.5)^2 = 1$$
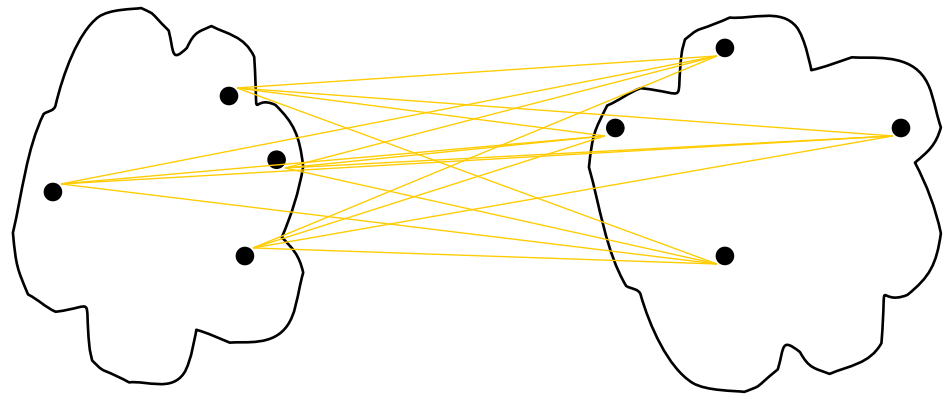
$$SSB = 2 \times (3-1.5)^2 + 2 \times (4.5-3)^2 = 9$$

$$Total = 1 + 9 = 10$$

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Unsupervised 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.
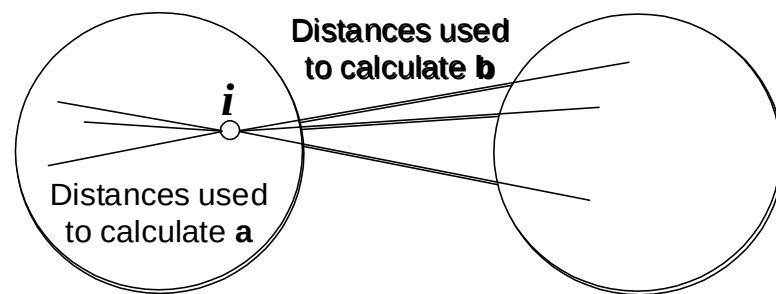
cohesion                                    separation

# Unsupervised Measures: Silhouette Coefficient

- Silhouette coefficient combines ideas of both cohesion and separation, but for individual points, as well as clusters and clusterings
- For an individual point, $i$
  - Calculate $a$ = average distance of $i$ to the points in its cluster
  - Calculate $b$ = min (average distance of $i$ to points in another cluster)
  - The silhouette coefficient for a point is then given by

    s = (b – a) / max(a,b)

  - Value can vary between -1 and 1
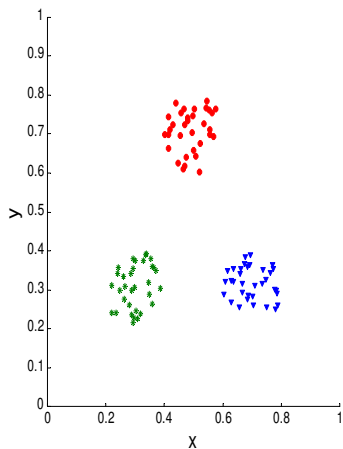  - Typically ranges between 0 and 1.
  - The closer to 1 the better.



Distances used to calculate **b**

$i$

Distances used to calculate **a**

- Can calculate the average silhouette coefficient for a cluster or a clustering

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Measuring Cluster Validity Via Correlation

- Two matrices
  - Proximity Matrix
  - Ideal Similarity 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 magnitude of correlation indicates that points that belong to the same cluster are close to each other.
  - Correlation may be positive or negative depending on whether the similarity matrix is a similarity or dissimilarity matrix

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

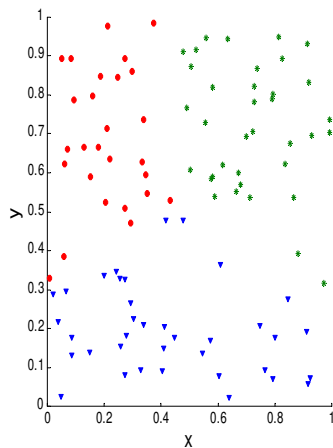# Measuring Cluster Validity via Correlation

- Correlation of ideal similarity and proximity matrices for the K-means clusterings of the following well-clustered data set.
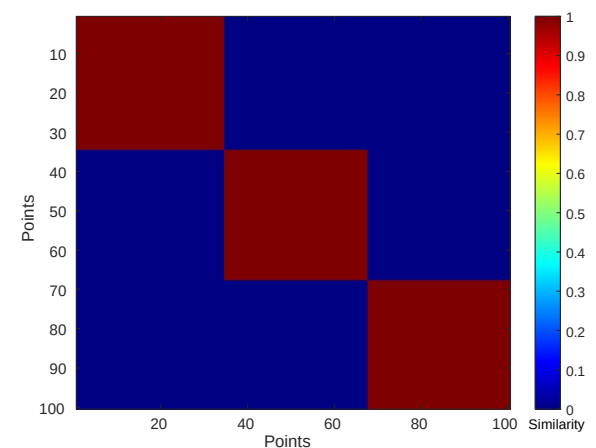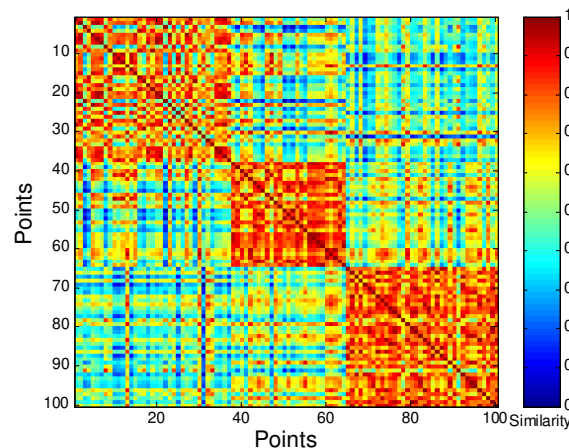


**Corr = 0.9235**

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Measuring Cluster Validity Via Correlation

- Correlation of ideal similarity and proximity matrices for the K-means clusterings of the following random data set.
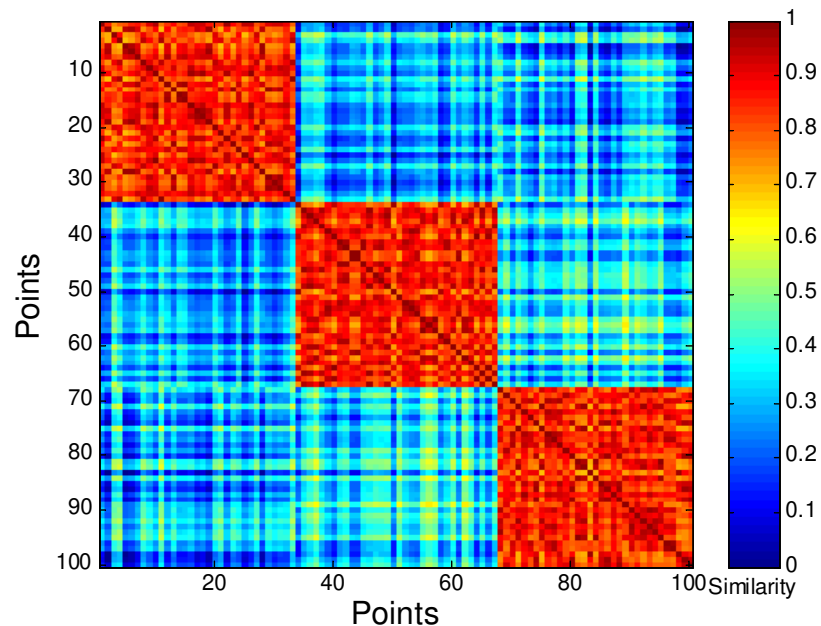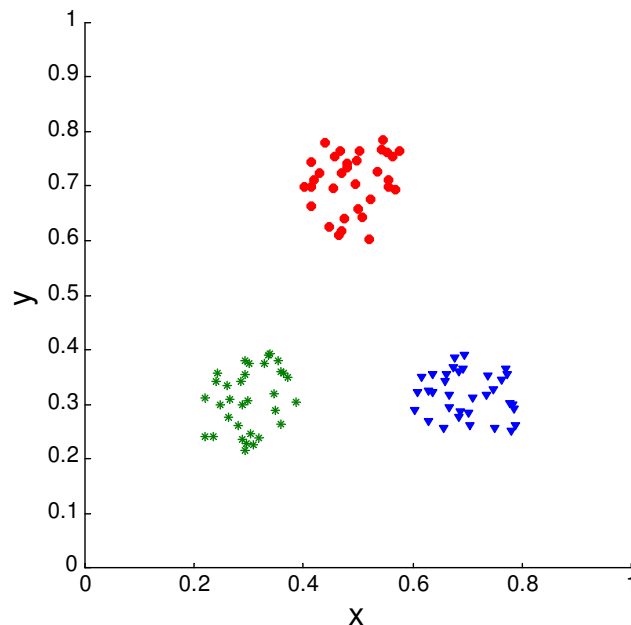


**K-means**

**Corr = 0.5810**

**Introduction to Data Mining, 2nd Edition**
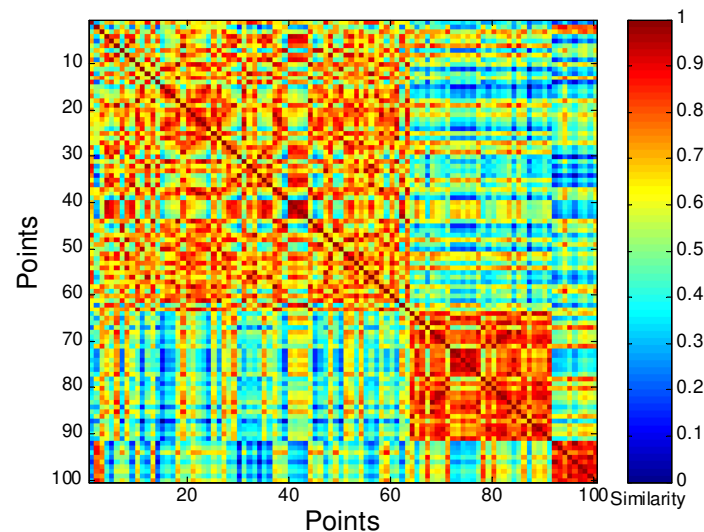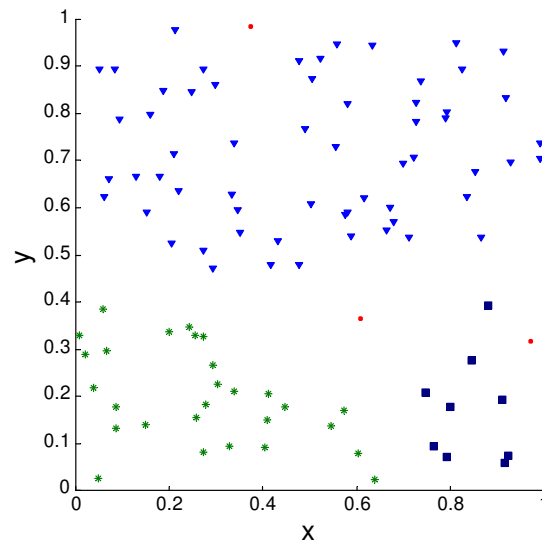**Tan, Steinbach, Karpatne, Kumar**

# Judging a Clustering Visually by its Similarity Matrix

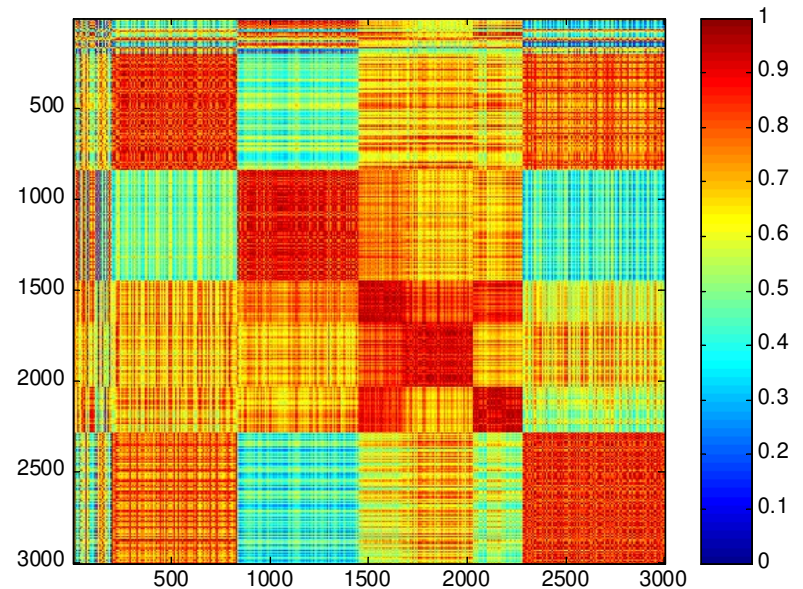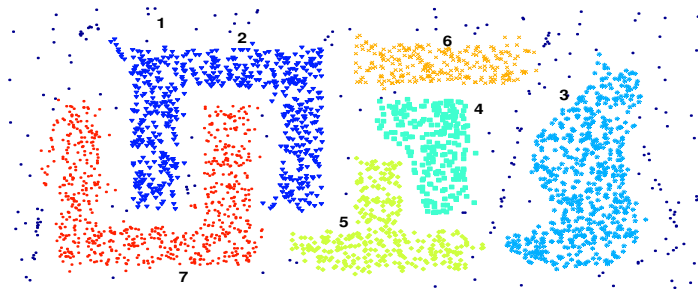- Order the similarity matrix with respect to cluster labels and inspect visually.

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

- Clusters in random data are not so crisp
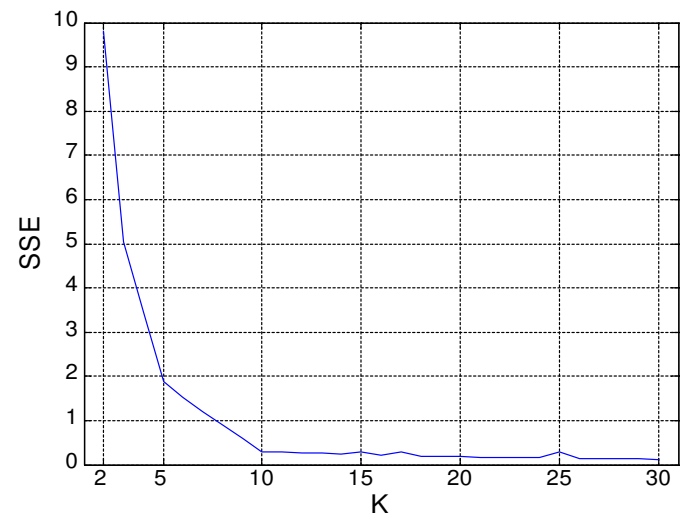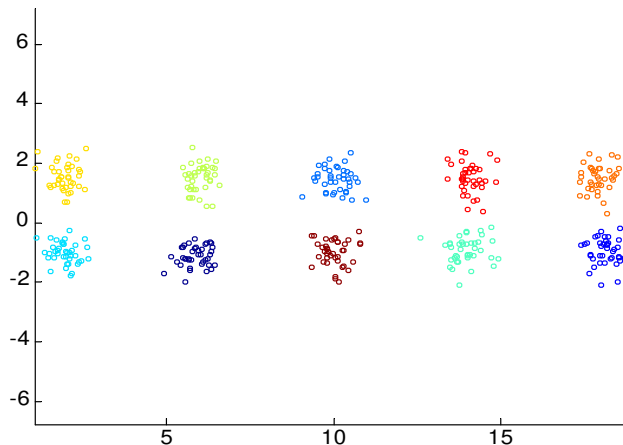


**DBSCAN**

# Judging a Clustering Visually by its Similarity Matrix
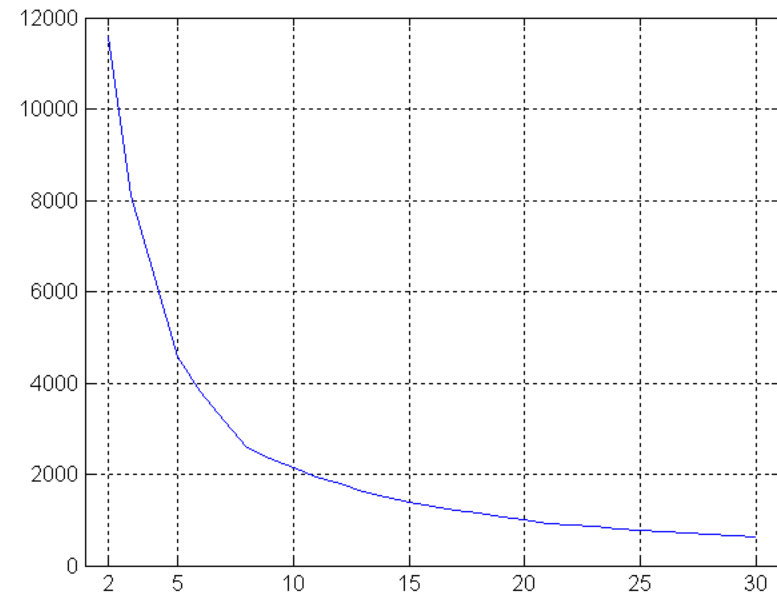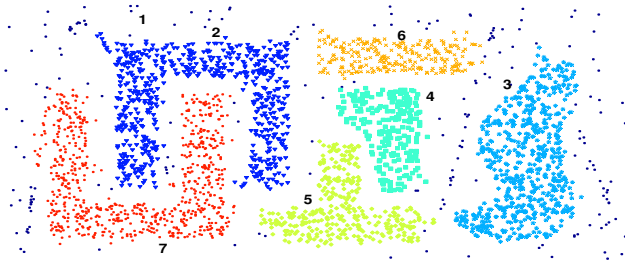


**DBSCAN**

# Determining the Correct Number of Clusters

- SSE is good for comparing two clusterings or two clusters
- SSE can also be used to estimate the number of clusters

# Determining the Correct Number of Clusters

- SSE curve for a more complicated data set



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

**Introduction to Data Mining, 2nd Edition**
**Tan, Steinbach, Karpatne, Kumar**

# Supervised Measures of Cluster Validity: Entropy and Purity

Table 5.9.  K-means Clustering Results for LA Document Data Set

| Cluster | Entertainment | Financial | Foreign | Metro | National | Sports | Entropy | Purity |
|---------|---------------|-----------|---------|-------|----------|--------|---------|--------|
| 1 | 3 | 5 | 40 | 506 | 96 | 27 | 1.2270 | 0.7474 |
| 2 | 4 | 7 | 280 | 29 | 39 | 2 | 1.1472 | 0.7756 |
| 3 | 1 | 1 | 1 | 7 | 4 | 671 | 0.1813 | 0.9796 |
| 4 | 10 | 162 | 3 | 119 | 73 | 2 | 1.7487 | 0.4390 |
| 5 | 331 | 22 | 5 | 70 | 13 | 23 | 1.3976 | 0.7134 |
| 6 | 5 | 358 | 12 | 212 | 48 | 13 | 1.5523 | 0.5525 |
| Total | 354 | 555 | 341 | 943 | 273 | 738 | 1.1450 | 0.7203 |

**entropy** For each cluster, the class distribution of the data is calculated first, i.e., for cluster $j$ we compute $p_{ij}$, the 'probability' that a member of cluster $j$ belongs to class $i$ as follows: $p_{ij} = m_{ij}/m_j$, where $m_j$ is the number of values in cluster $j$ and $m_{ij}$ is the number of values of class $i$ in cluster $j$. Then using this class distribution, the entropy of each cluster $j$ is calculated using the standard formula $e_j = \sum_{i=1}^{L} p_{ij} \log_2 p_{ij}$, where the $L$ is the number of classes. The total entropy for a set of clusters is calculated as the sum of the entropies of each cluster weighted by the size of each cluster, i.e., $e = \sum_{i=1}^{K} \frac{m_i}{m} e_j$, where $m_j$ is the size of cluster $j$, $K$ is the number of clusters, and $m$ is the total number of data points.

**purity** Using the terminology derived for entropy, the purity of cluster $j$, is given by $purity_j = \max p_{ij}$ and the overall purity of a clustering by $purity = \sum_{i=1}^{K} \frac{m_i}{m} purity_j$.

**Introduction to Data Mining, 2nd Edition
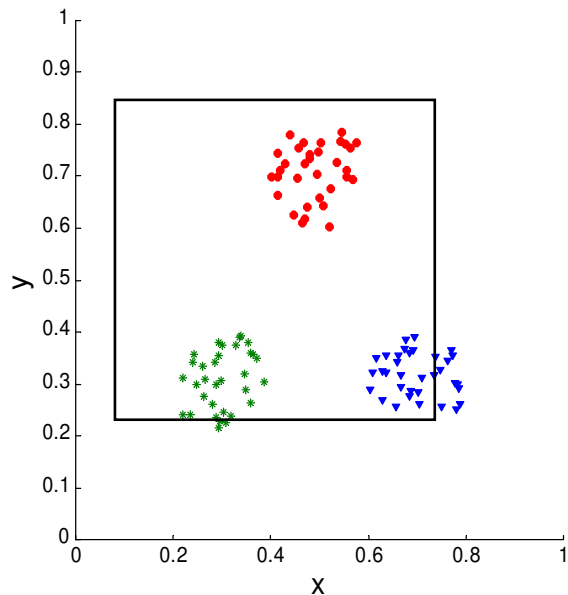Tan, Steinbach, Karpatne, Kumar**

# Assessing the Significance of Cluster Validity Measures

- 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
  - Compare the value of an index obtained from the given data with those resulting from random data.
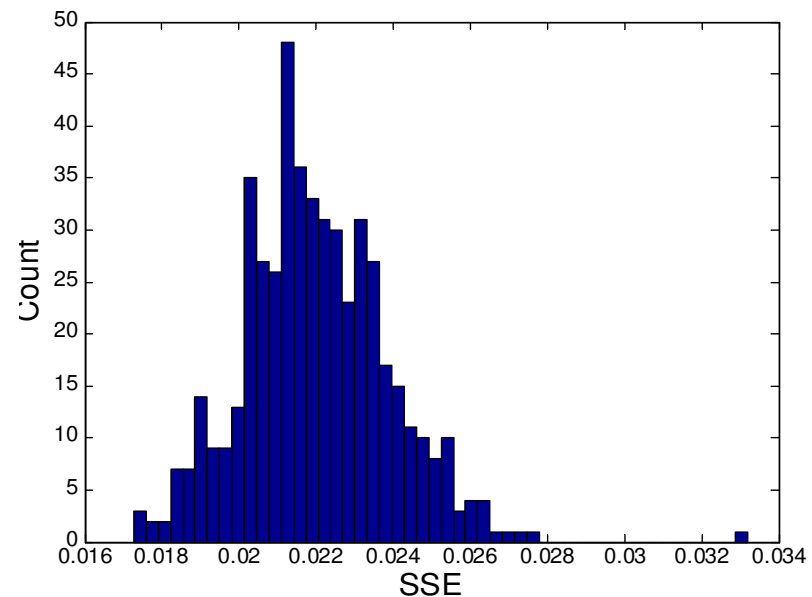    - If the value of the index is unlikely, then the cluster results are valid

# Statistical Framework for SSE

## Example

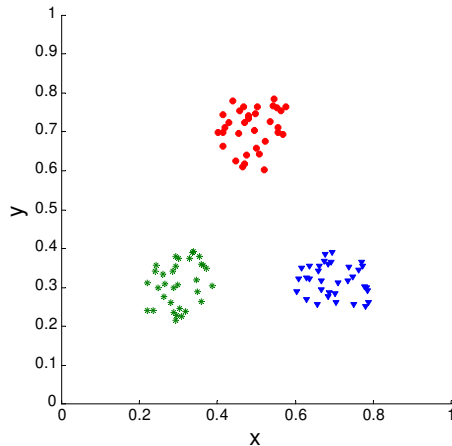- Compare SSE of three cohesive clusters against three clusters in random data



SSE = 0.005

Histogram shows SSE of three clusters in 500 sets of random data points of size 100 distributed over the range 0.2 – 0.8 for x and y values

**Introduction to Data Mining, 2nd Edition**
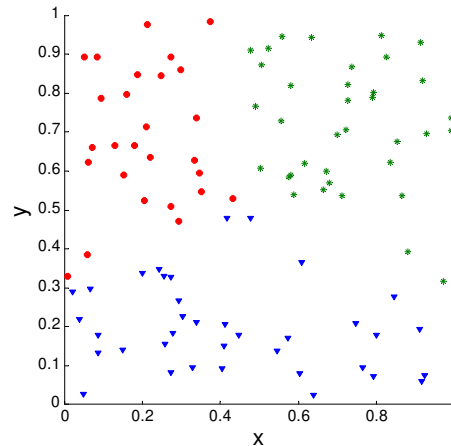**Tan, Steinbach, Karpatne, Kumar**

# Statistical Framework for Correlation

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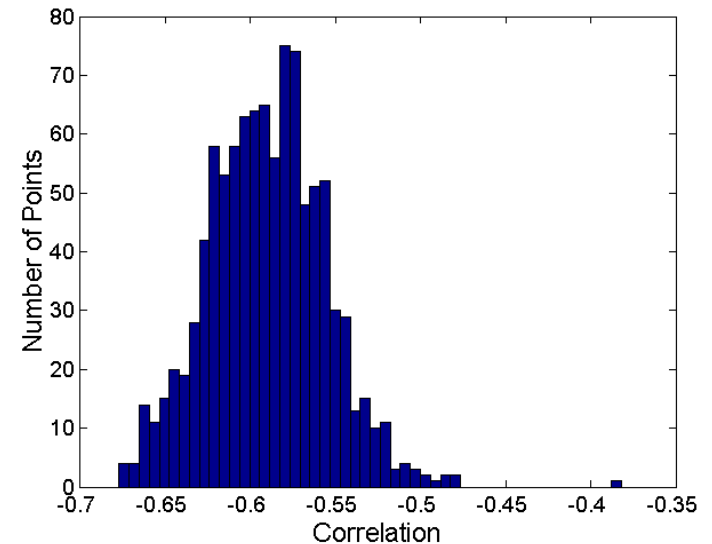


Histogram of correlation for 500 random data sets of size 100 with *x* and *y* values of points between 0.2 and 0.8.

**Corr = -0.9235**          **Corr = -0.5810**

Correlation is negative because it is calculated between a distance matrix and the ideal similarity matrix. Higher magnitude is better.

# 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**

- H. Xiong and Z. Li. *Clustering Validation Measures*. In C. C. Aggarwal and C. K. Reddy, editors, Data Clustering: Algorithms and Applications, pages 571–605. Chapman & Hall/CRC, 2013.