

# CSE 473: Pattern Recognition

# Unsupervised Learning:

## *Clustering*

# Proximity Measures Between Discrete-Valued Vectors

- Let  $F = \{0, 1, \dots, k-1\}$  be a set of symbols and  $X = \{\underline{x}_1, \dots, \underline{x}_N\} \subset F^l$
- Let  $A(\underline{x}, \underline{y}) = [a_{ij}]$ ,  $i, j = 0, 1, \dots, k-1$ , where  $a_{ij}$  is the number of places where  $\underline{x}$  has the  $i$ -th symbol and  $\underline{y}$  has the  $j$ -th symbol.

Example:  $l=6, k=3$

$$\begin{aligned} \mathbf{x} &= [0, 1, 2, 1, 2, 1]^T \\ \mathbf{y} &= [1, 0, 2, 1, 0, 1]^T \end{aligned} \quad A(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 2 & 0 \\ 1 & 0 & 1 \end{bmatrix}$$

NOTE: 
$$\sum_{i=0}^{k-1} \sum_{j=0}^{k-1} a_{ij} = l$$

# Proximity Measures Between Discrete-Valued Vectors

- Several proximity measures can be expressed as combinations of the elements of  $A(\underline{x}, \underline{y})$ .
  - Dissimilarity measures:
    - The **Hamming distance** (number of places where  $\underline{x}$  and  $\underline{y}$  differ)

$$d_H(\underline{x}, \underline{y}) = \sum_{i=0}^{k-1} \sum_{\substack{j=0 \\ j \neq i}}^{k-1} a_{ij}$$

# Proximity Measures Between Discrete-Valued Vectors

- Several proximity measures can be expressed as combinations of the elements of  $A(\underline{x}, \underline{y})$ .
  - Dissimilarity measures:
    - The **Hamming distance** (number of places where  $\underline{x}$  and  $\underline{y}$  differ)

$$d_H(\underline{x}, \underline{y}) = \sum_{i=0}^{k-1} \sum_{\substack{j=0 \\ j \neq i}}^{k-1} a_{ij}$$

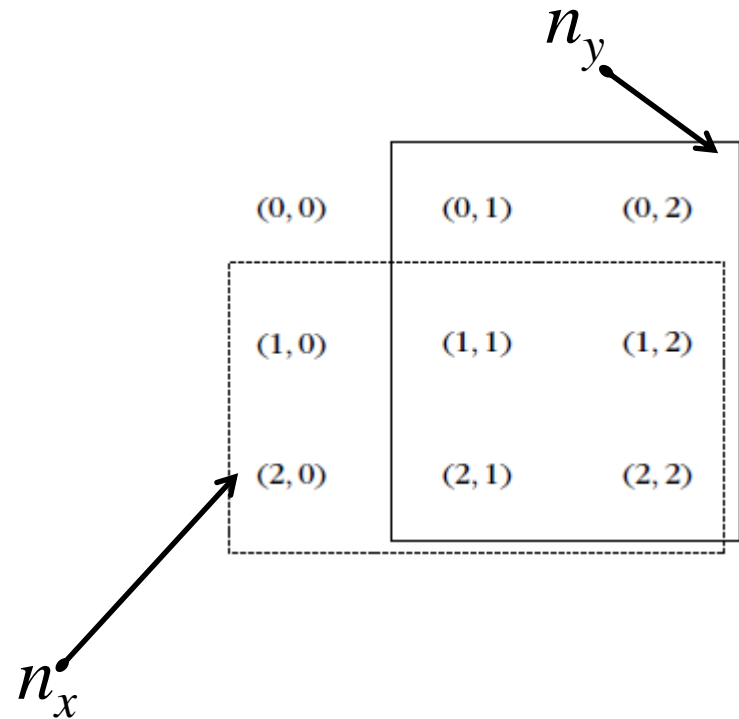
$$\begin{aligned}\mathbf{x} &= [0, 1, 2, 1, 2, 1]^T \\ \mathbf{y} &= [1, 0, 2, 1, 0, 1]^T\end{aligned}$$

$$A(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 2 & 0 \\ 1 & 0 & 1 \end{bmatrix}$$

# Proximity Measures Between Discrete-Valued Vectors

- Similarity measures:

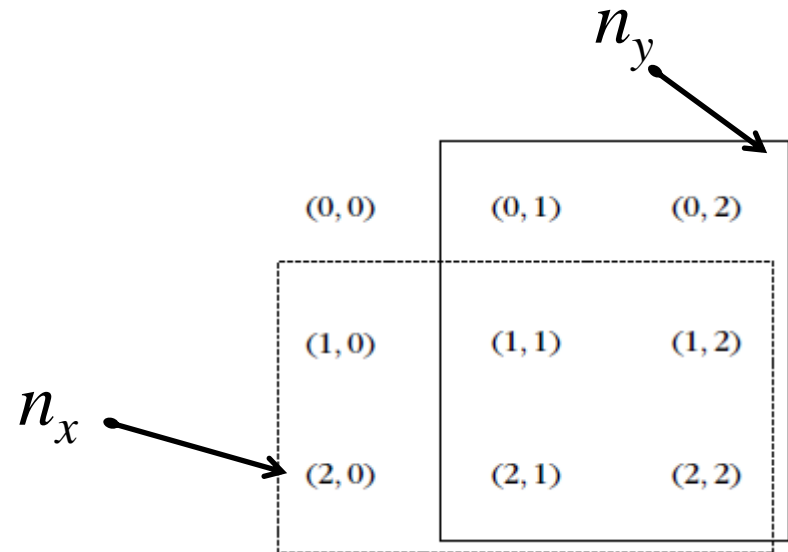
$$A(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 2 & 0 \\ 1 & 0 & 1 \end{bmatrix}$$



# Proximity Measures Between Discrete-Valued Vectors

- Similarity measures:

$$A(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 2 & 0 \\ 1 & 0 & 1 \end{bmatrix}$$



- Tanimoto measure :

$$s_T(\underline{x}, \underline{y}) = \frac{\sum_{i=1}^{k-1} a_{ii}}{n_x + n_y - \sum_{i=1}^{k-1} \sum_{j=1}^{k-1} a_{ij}}$$

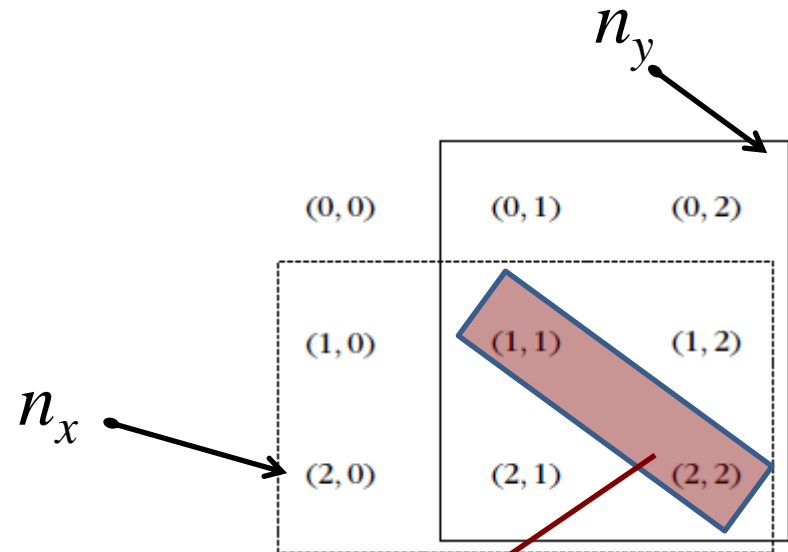
where

$$n_x = \sum_{i=1}^{k-1} \sum_{j=0}^{k-1} a_{ij}, \quad n_y = \sum_{i=0}^{k-1} \sum_{j=1}^{k-1} a_{ij},$$

# Proximity Measures Between Discrete-Valued Vectors

– Similarity measures:

$$A(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 2 & 0 \\ 1 & 0 & 1 \end{bmatrix}$$



• Tanimoto measure :

$$s_T(\underline{x}, \underline{y}) = \frac{\sum_{i=1}^{k-1} a_{ii}}{n_x + n_y - \sum_{i=1}^{k-1} \sum_{j=1}^{k-1} a_{ij}}$$

where

$$n_x = \sum_{i=1}^{k-1} \sum_{j=0}^{k-1} a_{ij}, \quad n_y = \sum_{i=0}^{k-1} \sum_{j=1}^{k-1} a_{ij},$$



# Proximity Measures Between Mixed-Valued Vectors

- Some of the coordinates of the vectors  $\underline{x}$  are real and the rest are discrete.

*Methods for measuring the proximity between two such  $\underline{x}_i$  and  $\underline{x}_j$ :*

- Adopt a proximity measure (PM) suitable for real-valued vectors.
- Convert the real-valued features to discrete ones and employ a discrete PM.

The more general case of mixed-valued vectors:

- Here nominal, ordinal, interval-scaled, ratio-scaled features are treated separately.

# Proximity Measures Between Mixed-Valued Vectors

The similarity function between  $\underline{x}_i$  and  $\underline{x}_j$  is:

$$s(\underline{x}_i, \underline{x}_j) = \sum_{q=1}^l s_q(\underline{x}_i, \underline{x}_j) / \sum_{q=1}^l w_q$$

In the above definition:

- $w_q=0$ , if one of the  $q$ -th coordinates of  $\underline{x}_i$  and  $\underline{x}_j$  are undefined or both the  $q$ -th coordinates are equal to 0. Otherwise  $w_q=1$ .
- If the  $q$ -th coordinates are **binary**,  $s_q(\underline{x}_i, \underline{x}_j)=1$  if  $x_{iq}=x_{jq}=1$  and 0 otherwise.
- If the  $q$ -th coordinates are **nominal or ordinal**,  $s_q(\underline{x}_i, \underline{x}_j)=1$  if  $x_{iq}=x_{jq}$  and 0 otherwise.
- If the  $q$ -th coordinates are interval or ratio scaled-valued

$$s_q(\underline{x}_i, \underline{x}_j) = 1 - |x_{iq} - x_{jq}| / r_q,$$

where  $r_q$  is the interval where the  $q$ -th coordinates of the vectors of the data set  $X$  lie.

# Fuzzy Proximity Measures

Let  $\underline{x}, \underline{y} \in [0, 1]^l$ . Here the value of the  $i$ -th coordinate,  $x_i$ , of  $\underline{x}$ , **is not the outcome of a measuring device**.

- The closer the coordinate  $x_i$  is to 1 (0), the more likely the vector  $\underline{x}$  **possesses** (does not possess) the  $i$ -th characteristic.
- As  $x_i$  approaches 0.5, the certainty about the possession or not of the  $i$ -th feature from  $\underline{x}$  decreases.

# Fuzzy Proximity Measures

A possible similarity measure that can quantify the above is:

$$s(x_i, y_i) = \max(\min(1 - x_i, 1 - y_i), \min(x_i, y_i))$$

Then

$$s_F^q(\underline{x}, \underline{y}) = \left( \sum_{i=1}^l s(x_i, y_i)^q \right)^{1/q}$$

# Proximity Measures For Missing Data

- For some vectors of the data set  $X$ , some features values are unknown

*Ways to face the problem:*

- Discard all vectors with missing values (not recommended for small data sets)
- Find the mean value  $m_i$  of the available  $i$ -th feature values over that data set and substitute the missing  $i$ -th feature values with  $m_i$ .

# Proximity Measures For Missing Data

- Define  $b_i=0$ , if both the  $i$ -th features  $x_i, y_i$  are available and  $b_i=1$  otherwise. Then

$$\wp(\underline{x}, \underline{y}) = \frac{l}{l - \sum_{i=1}^l b_i} \sum_{\text{all } i: b_i=0} \phi(x_i, y_i)$$

where  $\phi(x_i, y_i)$  denotes the PM between two scalars  $x_i, y_i$ .

# Proximity Measures For Missing Data

- Find the average proximities  $\phi_{avg}(i)$  between all feature vectors in  $X$  along all components. Then

$$\wp(\underline{x}, \underline{y}) = \sum_{i=1}^l \psi(x_i, y_i)$$

where  $\psi(x_i, y_i) = \phi(x_i, y_i)$ , if both  $x_i$  and  $y_i$  are available and  $\phi_{avg}(i)$  otherwise.

# Proximity Functions Between A Vector and A Set

- Let  $X = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$  and  $C \subset X$ ,  $\underline{x} \in X$
- All points of  $C$  contribute to the definition of  $\wp(\underline{x}, C)$

– **Max** proximity function

$$\wp_{\max}^{ps}(\underline{x}, C) = \max_{\underline{y} \in C} \wp(\underline{x}, \underline{y})$$

– **Min** proximity function

$$\wp_{\min}^{ps}(\underline{x}, C) = \min_{\underline{y} \in C} \wp(\underline{x}, \underline{y})$$

– **Average** proximity function

$$\wp_{avg}^{ps}(\underline{x}, C) = \frac{1}{n_C} \sum_{\underline{y} \in C} \wp(\underline{x}, \underline{y}) \quad (n_C \text{ is the cardinality of } C)$$



- A representative(s) of  $C$ ,  $r_C$ , contributes to the definition of  $\wp(\underline{x}, C)$

In this case:  $\wp(\underline{x}, C) = \wp(\underline{x}, r_C)$

Typical representatives are:

- The mean vector:

$$\underline{m}_p = \left( \frac{1}{n_C} \right) \sum_{y \in C} \underline{y}$$

where  $n_C$  is the cardinality of  $C$

- The mean center:

$$\underline{m}_C \in C : \sum_{y \in C} d(\underline{m}_C, \underline{y}) \leq \sum_{y \in C} d(\underline{z}, \underline{y}), \quad \forall \underline{z} \in C$$

$d$ : a dissimilarity measure

- The median center:

$$\underline{m}_{med} \in C : med(d(\underline{m}_{med}, \underline{y}) \mid \underline{y} \in C) \leq med(d(\underline{z}, \underline{y}) \mid \underline{y} \in C), \quad \forall \underline{z} \in C$$

**NOTE:** Other representatives (e.g., hyperplanes, hyperspheres) are useful in certain applications (e.g., object identification using clustering techniques).

# Proximity Functions Between Sets

- Let  $X = \{\underline{x}_1, \dots, \underline{x}_N\}$ ,  $D_i, D_j \subset X$  and  $n_i = |D_i|$ ,  $n_j = |D_j|$
- All points of each set contribute to  $\wp(D_i, D_j)$ 
  - **Max** proximity function

$$\wp_{\max}^{ss}(D_i, D_j) = \max_{\underline{x} \in D_i, \underline{y} \in D_j} \wp(\underline{x}, \underline{y})$$

- **Min** proximity function

$$\wp_{\min}^{ss}(D_i, D_j) = \min_{\underline{x} \in D_i, \underline{y} \in D_j} \wp(\underline{x}, \underline{y})$$

- **Average** proximity function

$$\wp_{\text{avg}}^{ss}(D_i, D_j) = \left( \frac{1}{n_i n_j} \right) \sum_{\underline{x} \in D_i} \sum_{\underline{y} \in D_j} \wp(\underline{x}, \underline{y})$$

# Proximity Functions Between Sets

- Each set  $D_i$  is represented by its representative vector  $\underline{m}_i$ 
  - Mean proximity function (it is a measure provided that  $\wp$  is a measure):

$$\wp_{mean}^{ss}(D_i, D_j) = \wp(\underline{m}_i, \underline{m}_j)$$

## ➤ Remarks:

- Different choices of proximity functions between sets may lead to totally different clustering results.
- Different proximity measures between vectors in the same proximity function between sets may lead to totally different clustering results.
- The only way to achieve a proper clustering is
  - by trial and error and,
  - taking into account the opinion of an expert in the field of application.

# CLUSTERING ALGORITHMS

- Number of possible clusterings

Let  $X = \{x_1, x_2, \dots, x_N\}$ .

**Question:** In how many ways the  $N$  points can be assigned into  $m$  groups?

**Answer: ?**

# CLUSTERING ALGORITHMS

- Number of possible clusterings

Let  $X = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$ .

**Question:** In how many ways the  $N$  points can be assigned into  $m$  groups?

**Answer: ?**

Let  $S(N, m)$  = number of possible ways of clusterings

$L_N^m$  is the list of all possible clusterings of  $N$  vectors into  $m$  groups

# CLUSTERING ALGORITHMS

- Number of possible clusterings

Let  $X = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$ .

**Question:** In how many ways the  $N$  points can be assigned into  $m$  groups?

**Answer: ?**

Let  $S(N, m)$  = number of possible ways of clusterings

$L_N^m$  is the list of all possible clusterings of  $N$  vectors into  $m$  groups

Therefore,

$L_{N-1}^{m-1}$  is the list containing all possible clusterings of the  $N - 1$  vectors into  $m - 1$  groups

# CLUSTERING ALGORITHMS

- Number of possible clusterings

Let  $X = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$ .

**Question:** In how many ways the  $N$  points can be assigned into  $m$  groups?

**Answer: ?**

$L_N^m$  is the list of all possible clusterings of  $N$  vectors into  $m$  groups

Therefore,

$L_{N-1}^{m-1}$  be the list containing all possible clusterings of the  $N-1$  vectors into  $m-1$  groups

$L_{N-1}^m$  be the list containing all possible clusterings of the  $N-1$  vectors into  $m$  groups



# CLUSTERING ALGORITHMS

- Number of possible clusterings

Let  $X = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$ .

**Question:** In how many ways the  $N$  points can be assigned into  $m$  groups?

**Answer: ?**

Now find  $L_N^m$  from  $L_{N-1}^{m-1}$  and  $L_{N-1}^m$ , in two ways:

- Either will be added to one of the clusters of any member of  $L_{N-1}^m$
- Or will form a new cluster to each member of  $L_{N-1}^{m-1}$

# CLUSTERING ALGORITHMS

- Number of possible clusterings

Let  $X = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$ .

**Question:** In how many ways the  $N$  points can be assigned into  $m$  groups?

**Answer: ?**

Now find  $L_N^m$  from  $L_{N-1}^{m-1}$  and  $L_{N-1}^m$ , in two ways:

- Either will be added to one of the clusters of any member of  $L_{N-1}^m$
- Or will form a new cluster to each member of  $L_{N-1}^{m-1}$

Therefore,  $S(N, m) = mS(N-1, m) + S(N-1, m-1)$

# CLUSTERING ALGORITHMS

- Number of possible clusterings

Let  $X = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$ .

**Question:** In how many ways the  $N$  points can be assigned into  $m$  groups?

**Answer:**

$$S(N, m) = \frac{1}{m!} \sum_{i=0}^m (-1)^{m-i} \binom{m}{i} i^N$$

# CLUSTERING ALGORITHMS

- Number of possible clusterings

Let  $X = \{x_1, x_2, \dots, x_N\}$ .

**Question:** In how many ways the  $N$  points can be assigned into  $m$  groups?

**Answer:**

$$S(N, m) = \frac{1}{m!} \sum_{i=0}^m (-1)^{m-i} \binom{m}{i} i^N$$

— Examples:

$$S(15, 3) = 2,375,101$$

# CLUSTERING ALGORITHMS

- Number of possible clusterings

Let  $X = \{x_1, x_2, \dots, x_N\}$ .

**Question:** In how many ways the  $N$  points can be assigned into  $m$  groups?

**Answer:**

$$S(N, m) = \frac{1}{m!} \sum_{i=0}^m (-1)^{m-i} \binom{m}{i} i^N$$

— Examples:

$$S(15, 3) = 2,375,101$$

$$S(20, 4) = 45,232,115,901$$

# CLUSTERING ALGORITHMS

- Number of possible clusterings

Let  $X = \{x_1, x_2, \dots, x_N\}$ .

**Question:** In how many ways the  $N$  points can be assigned into  $m$  groups?

**Answer:**

$$S(N, m) = \frac{1}{m!} \sum_{i=0}^m (-1)^{m-i} \binom{m}{i} i^N$$

— Examples:

$$S(15, 3) = 2,375,101$$

$$S(20, 4) = 45,232,115,901$$

$$S(100, 5) = 10^{68} !!$$

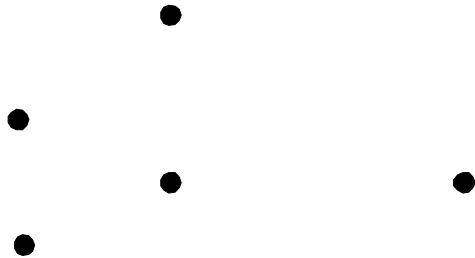
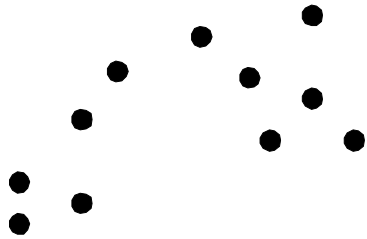
- A way out:
  - Consider only a small fraction of clusterings of  $X$  and select a “sensible” clustering among them.
    - Question 1: Which fraction of clusterings is considered?
    - Question 2: What “sensible” means?
  - The answer depends on the specific clustering algorithm and the specific criteria to be adopted.

# Clustering Algorithms

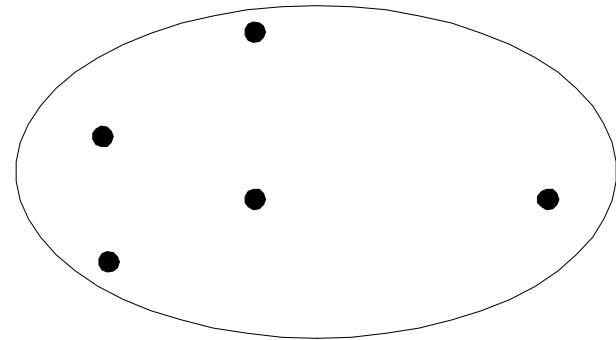
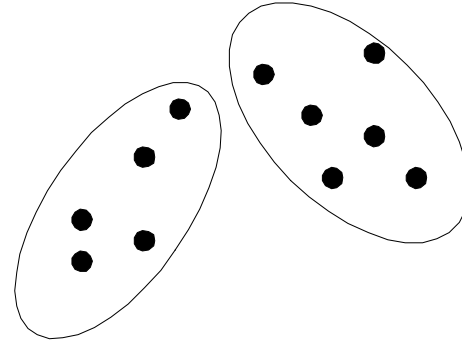
- 2 Important distinctions: **hierarchical** and **partitional**
- Partitional Clustering
  - A division data objects into **non-overlapping subsets** (clusters) such that each data object is in exactly one subset
- Hierarchical clustering
  - A set of **nested clusters** organized as a hierarchical tree



# Partitional Clustering

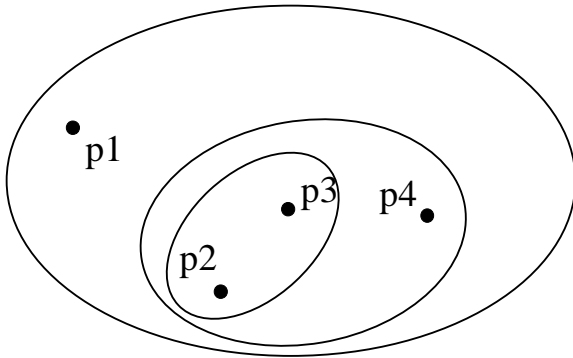


Original Points

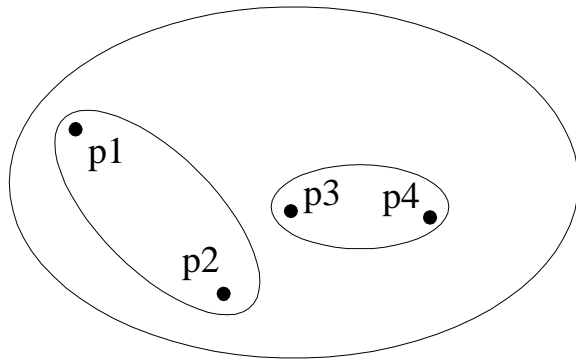


A Partitional Clustering

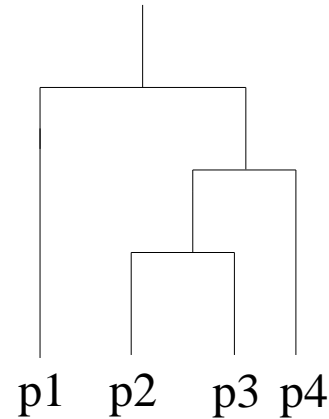
# Hierarchical Clustering



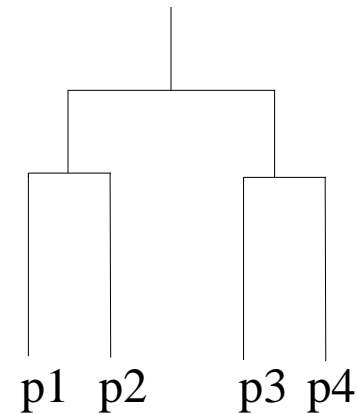
Traditional Hierarchical Clustering



Non-traditional Hierarchical Clustering



Traditional **Dendrogram**



Non-traditional **Dendrogram**

# Hierarchical Clustering

- **Hierarchical:** A sequence of (nested) clusterings is produced.
  - Agglomerative
  - Divisive
  - Combinations of the above (e.g., the Chameleon algorithm.)

## OTHER MAJOR CATEGORIES OF CLUSTERING ALGORITHMS

- **Sequential:** A single clustering is produced. One or few sequential passes on the data.

- **Cost function optimization.** For most of the cases a *single* clustering is obtained.
  - **Hard clustering** (each point belongs exclusively to a single cluster):
    - Basic hard clustering algorithms (e.g.,  $k$ -means)
    - Branch and bound
    - Boundary detection
    - Genetic clustering algorithms
    - ...
    - ..
    - .
  - **Fuzzy clustering** (each point belongs to more than one clusters simultaneously).
  - **Probabilistic clustering** (it is based on the *probability* of a point to belong to a cluster).

- Other schemes:
  - Exclusive versus non-exclusive
    - In non-exclusive clusterings, points may belong to multiple clusters.
    - Can represent multiple classes or ‘border’ points
  - Partial versus complete
    - In some cases, we only want to cluster some of the data
  - Heterogeneous versus homogeneous
    - Cluster of widely different sizes, shapes, and densities

- Other schemes:
  - Algorithms based on graph theory (e.g., Minimum Spanning Tree, regions of influence, directed trees).
  - Density based clustering algorithms.
  - Subspace clustering algorithms.
  - Kernel based methods

# Clustering and Clusters :Difference?

## Types of Clusters

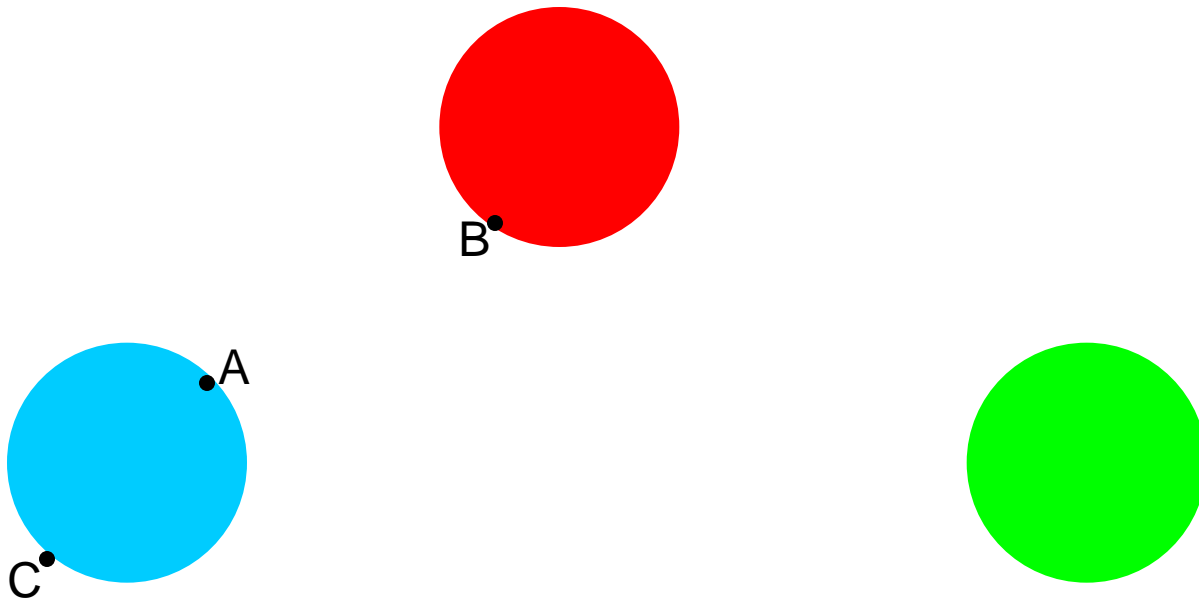
- Well-separated clusters
- Center-based clusters
- Contiguous clusters
- Density-based clusters
- Property or Conceptual
- Described by an Objective Function



# 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.

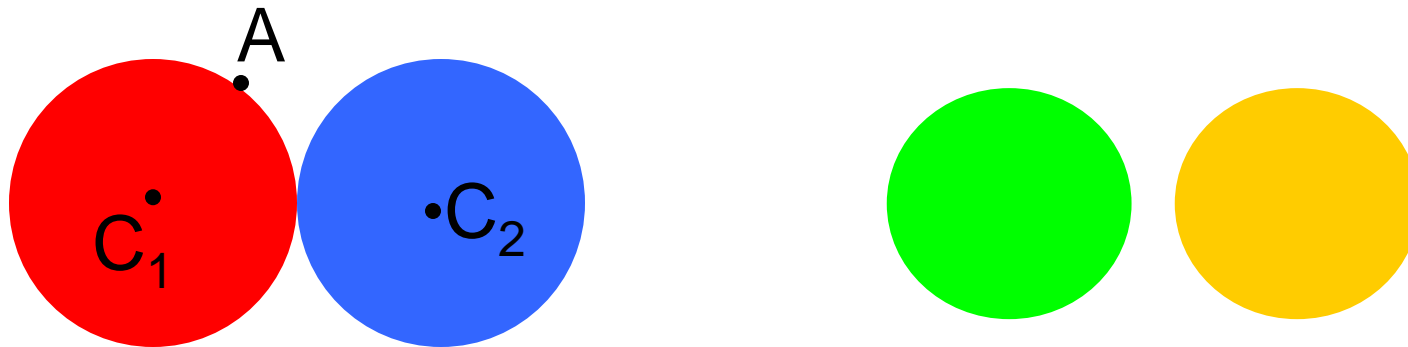


3 well-separated clusters

# Types of Clusters: Center-Based

- Center-based

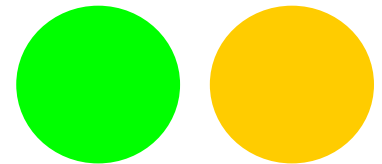
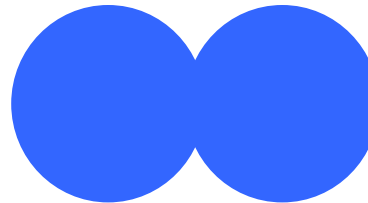
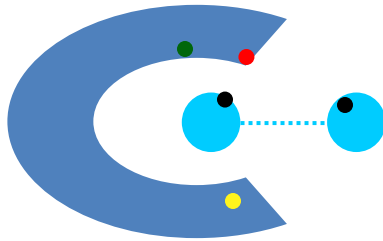
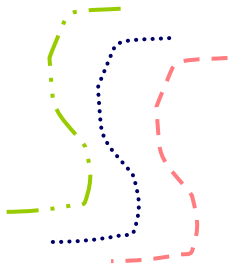
- A cluster is a set of objects such that an object in a cluster is closer (more similar) to the “center” of a cluster, than to the center of any other cluster
- The center of a cluster is often a centroid, the average of all the points in the cluster, or a medoid, the most “representative” point of a cluster



4 center-based clusters

# 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.

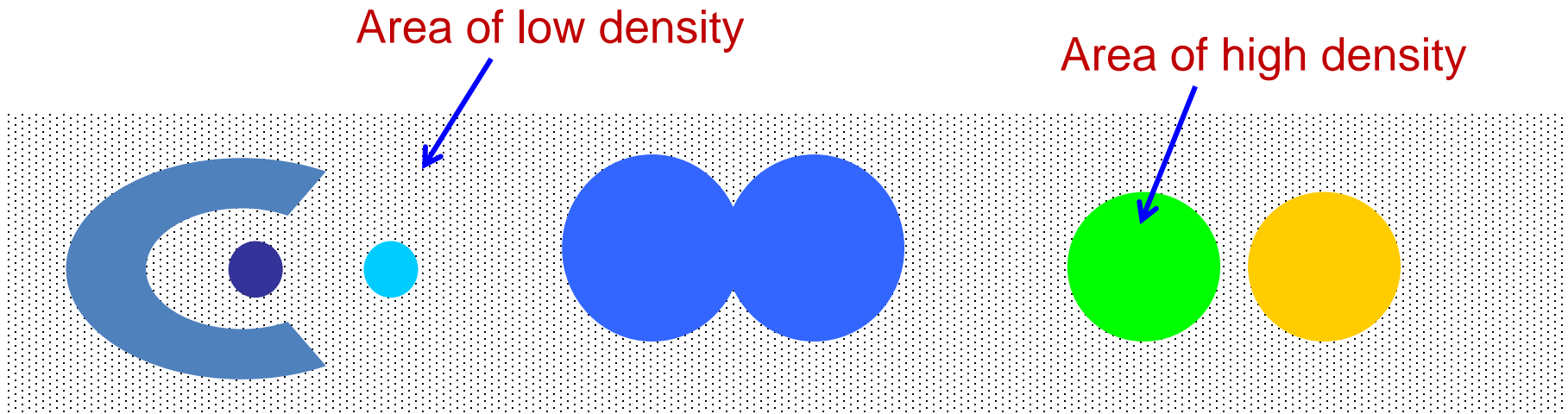


8 contiguous clusters

# Types of Clusters: Density-Based

- Density-based

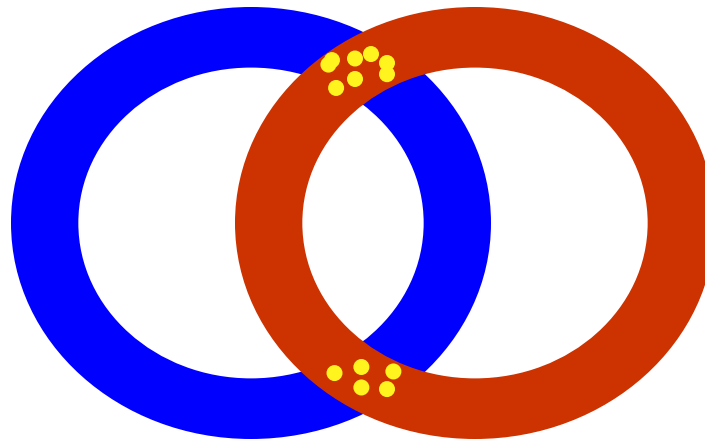
- A cluster is a **dense region of points**, which is **separated by low-density regions**, from other regions of high density.
- Used when the clusters are irregular or intertwined, and when noise and outliers are present.



6 density-based clusters

# Types of Clusters: Conceptual Clusters

- Shared Property or Conceptual Clusters
  - Finds clusters that share some **common property** or represent a particular concept.



2 Overlapping Circles

# SEQUENTIAL CLUSTERING ALGORITHMS

- Characteristics:
  - One or very few passes on the data
  - No a-prior knowledge about the number of clusters
  - The clusters are defined with the aid of
    - An appropriately defined distance  $d(\underline{x}, C)$  of a point from a cluster.
    - A threshold  $\theta$  associated with the distance.

## ➤ Basic Sequential Clustering Algorithm (BSAS)

- $m=1 \setminus \{\text{number of clusters}\}$
- $C_m = \{x_1\}$
- For  $i=2$  to  $N$ 
  - Find  $C_k$ :  $d(x_i, C_k) = \min_{1 \leq j \leq m} d(x_i, C_j)$
  - If  $(d(x_i, C_k) > \Theta)$  AND  $(m < q)$  then
    - o  $m = m + 1$
    - o  $C_m = \{x_i\}$
  - Else
    - o  $C_k = C_k \cup \{x_i\}$
    - o Where necessary, update representatives (\*)
  - End {if}
- End {for}

## ➤ Basic Sequential Clustering Algorithm (BSAS)

- $m=1$  \{number of clusters\}
- $C_m = \{\underline{x}_1\}$
- For  $i=2$  to  $N$ 
  - Find  $C_k$ :  $d(\underline{x}_i, C_k) = \min_{1 \leq j \leq m} d(\underline{x}_i, C_j)$
  - If  $(d(\underline{x}_i, C_k) > \Theta)$  AND  $(m < q)$  then
    - o  $m = m + 1$
    - o  $C_m = \{\underline{x}_i\}$
  - Else
    - o  $C_k = C_k \cup \{\underline{x}_i\}$
    - o Where necessary, update representatives (\*)
  - End {if}
- End {for}

---

(\*) When the mean vector  $\underline{m}_C$  is used as representative of the cluster  $C$  with  $n_c$  elements, the updating in the light of a new vector  $\underline{x}$  becomes

$$\underline{m}_C^{new} = (n_C \underline{m}_C + \underline{x}) / (n_C + 1)$$



## ➤ Remarks:

- The **order of presentation of the data** in the algorithm plays important role in the clustering results.
  - Different order of presentation may lead to totally different clustering results

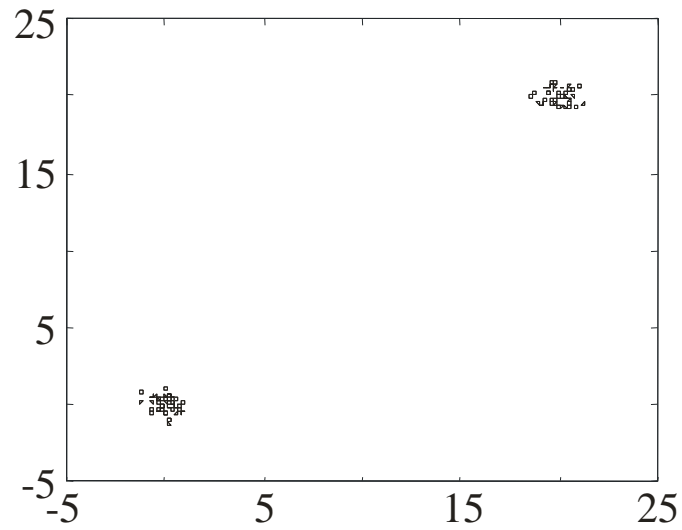
## ➤ Remarks:

- The **order of presentation of the data** in the algorithm plays important role in the clustering results.
  - Different order of presentation may lead to totally different clustering results
- In BSAS the decision for a vector  $\underline{x}$  is reached prior to the final cluster formation.

## ➤ Remarks:

- The **order of presentation of the data** in the algorithm plays important role in the clustering results.
  - Different order of presentation may lead to totally different clustering results
- In BSAS the decision for a vector  $\underline{x}$  is reached prior to the final cluster formation.
- BSAS perform a single pass on the data. Its complexity is  $O(N)$ .

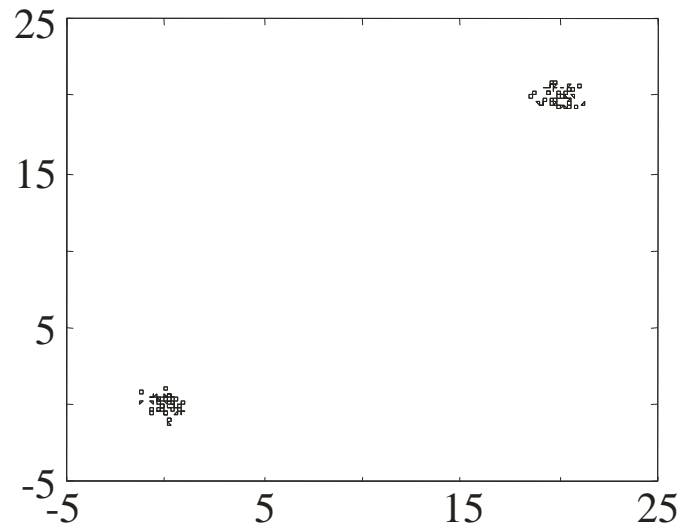
➤ Estimating the number of clusters in the data set:



➤ Estimating the number of clusters in the data set:

Let  $BSAS(\Theta)$  denote the  $BSAS$  algorithm when the dissimilarity threshold is  $\Theta$ .

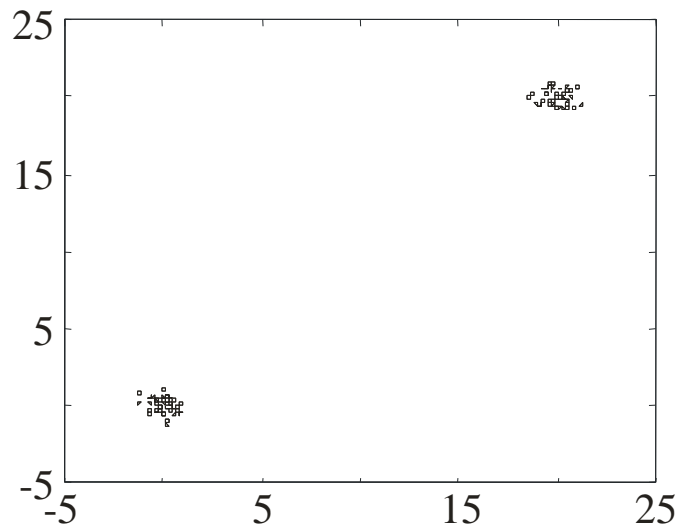
- Find
  - $a$ , the maximum and
  - $b$ , the minimum distances among the points



➤ Estimating the number of clusters in the data set:

Let  $BSAS(\Theta)$  denote the  $BSAS$  algorithm when the dissimilarity threshold is  $\Theta$ .

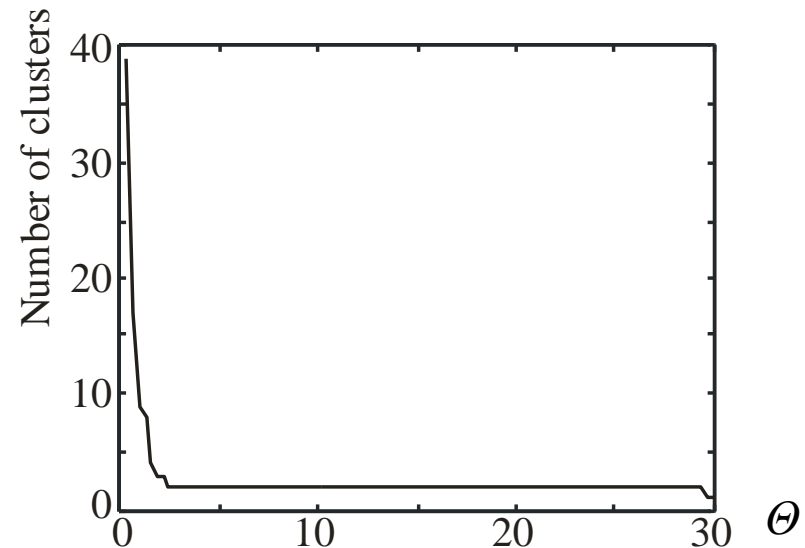
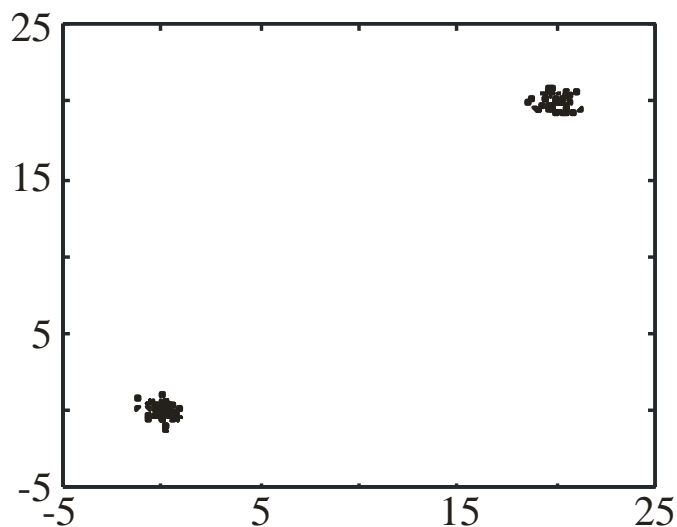
- For  $\Theta=a$  to  $b$  step  $c$ 
  - Run  $s$  times  $BSAS(\Theta)$ , each time presenting the data in a different order.
  - Estimate the number of clusters  $m_{\Theta}$ , as the most frequent number resulting from the  $s$  runs of  $BSAS(\Theta)$ .
- Next  $\Theta$



➤ Estimating the number of clusters in the data set:

Let  $BSAS(\Theta)$  denote the  $BSAS$  algorithm when the dissimilarity threshold is  $\Theta$ .

- For  $\Theta=a$  to  $b$  step  $c$ 
  - Run  $s$  times  $BSAS(\Theta)$ , each time presenting the data in a different order.
  - Estimate the number of clusters  $m_\Theta$ , as the most frequent number resulting from the  $s$  runs of  $BSAS(\Theta)$ .
- Next  $\Theta$
- Plot  $m_\Theta$  versus  $\Theta$  and identify the number of clusters  $m$  as the one corresponding to the widest flat region in the above graph.



## ➤ MBSAS, a modification of BSAS

In BSAS a decision for a data vector  $\underline{x}$  is reached prior to the final cluster formation, which is determined after all vectors have been presented to the algorithm.



## ➤ MBSAS, a modification of BSAS

In BSAS a decision for a data vector  $\underline{x}$  is reached prior to the final cluster formation, which is determined after all vectors have been presented to the algorithm.

- MBSAS deals with the above drawback, at the cost of presenting the data twice to the algorithm.

## ➤ MBSAS, a modification of BSAS

### MBSAS consists of:

- A **cluster determination phase** (first pass on the data), which is the same as BSAS with the exception that no vector is assigned to an already formed cluster. At the end of this phase, each cluster consists of a single element.
- A **pattern ~~classification~~ phase** (second pass on the data), where each one of the unassigned vector is assigned to its closest cluster.

## ➤ MBSAS, a modification of BSAS

- *Cluster Determination*

- $m=1 \setminus \{\text{number of clusters}\}$
- $C_m = \{\underline{x}_1\}$
- For  $i=2$  to  $N$ 
  - Find  $C_k$ :  $d(\underline{x}_i, C_k) = \min_{1 \leq j \leq m} d(\underline{x}_i, C_j)$
  - If  $(d(\underline{x}_i, C_k) > \Theta)$  AND  $(m < q)$  then
    - o  $m = m + 1$
    - o  $C_m = \{\underline{x}_i\}$
  - End {if}
- End {for}

## ➤ MBSAS, a modification of BSAS

- *Pattern ~~Classification~~ Assignment*
  - For  $i = 1$  to  $N$ 
    - If  $x_i$  has not been assigned to a cluster, then
      - Find  $C_k$ :  $d(x_i, C_k) = \min_{1 \leq j \leq m} d(x_i, C_j)$
      - $C_k = C_k \cup \{x_i\}$
      - where necessary, update representatives
    - End {if}
  - End {for}

## ➤ MBSAS, a modification of BSAS

### ➤ Remarks:

- decision for a vector  $\underline{x}$  during the pattern classification phase is reached taking into account all clusters.
- still sensitive to the order of presentation of the vectors.
- requires two passes on the data. complexity is  $O(N)$ .

## ➤ The maxmin algorithm, a variant of MBSAS

Let  $X$  = the set of all data points

$W$  = the set of all points that have been chosen to form clusters up to the current iteration step.

## ➤ The maxmin algorithm, a variant of MBSAS

### cluster formation:

- For each  $\underline{x} \in X - W$  determine  $d_x = \min_{\underline{z} \in W} d(\underline{x}, \underline{z})$
- Determine  $\underline{y}$ :  $d_y = \max_{\underline{x} \in X - W} d_x$
- If  $d_y$  is greater than a pre-specified threshold then
  - this vector forms a new cluster
- else
  - the cluster determination phase of the algorithm terminates.
- End {if}

### pattern classification:

assign each unassigned vector to its closest cluster

➤ Remarks:

- The maxmin algorithm is more computationally demanding than MBSAS.
- However, it is expected to produce better clustering results.



- A two-threshold sequential scheme (TTSAS)
  - The **formation** of the clusters, as well as the **assignment** of vectors to clusters, is carried out **concurrently** (like BSAS and unlike MBSAS)
  - Two thresholds  $\Theta_1$  and  $\Theta_2$  ( $\Theta_1 < \Theta_2$ ) are employed

- A two-threshold sequential scheme (TTSAS)
  - The **general idea** is the following:
    - If the distance  $d(\underline{x}, C)$  of  $\underline{x}$  from its closest cluster,  $C$ , is greater than  $\Theta_2$  then:
      - A new cluster represented by  $\underline{x}$  is formed.
    - Else if  $d(\underline{x}, C) < \Theta_1$  then
      - $\underline{x}$  is assigned to  $C$ .
    - Else
      - The decision is postponed to a later stage.
    - End {if}

- A two-threshold sequential scheme (TTSAS)
  - The **general idea** is the following:
    - If the distance  $d(\underline{x}, C)$  of  $\underline{x}$  from its closest cluster,  $C$ , is greater than  $\Theta_2$  then:
      - A new cluster represented by  $\underline{x}$  is formed.
    - Else if  $d(\underline{x}, C) < \Theta_1$  then
      - $\underline{x}$  is assigned to  $C$ .
    - Else
      - The decision is postponed to a later stage.
    - End {if}

The unassigned vectors are presented iteratively to the algorithm until all of them are classified.

➤ Remarks:

- In practice, a few passes ( $\geq 2$ ) of the data set are required.
- TTSAS is less sensitive to the order of data presentation, compared to BSAS.

- Refinement stages

Why necessary?

The problem of **closeness of clusters**: *“In all the above algorithms it may happen that two formed clusters lie very close to each other”.*

- Refinement stages

- A simple merging procedure

- (A) Find  $C_i, C_j$  ( $i < j$ ) such that  $d(C_i, C_j) = \min_{k,r=1,\dots,m, k \neq r} d(C_k, C_r)$
    - If  $d(C_i, C_j) \leq M_1$  then  $\{ M_1 \text{ is a user-defined threshold} \}$ 
      - Merge  $C_i, C_j$  to  $C_i$  and eliminate  $C_j$ .
      - If necessary, update the cluster representative of  $C_i$ .
      - Rename the clusters  $C_{j+1}, \dots, C_m$  to  $C_j, \dots, C_{m-1}$ , respectively.
      - $m = m - 1$
      - Go to (A)
    - Else
      - Stop
    - End {if}

- Reassignment of vectors

Why necessary?

The problem of sensitivity to the order of data presentation:

*“A vector  $\underline{x}$  may have been assigned to a cluster  $C_i$  at the current stage but another cluster  $C_j$  may be formed at a later stage that lies closer to  $\underline{x}$ ”*

– A simple reassignment procedure

- For  $i=1$  to  $N$ 
  - Find  $C_j$  such that  $d(\underline{x}_i, C_j) = \min_{k=1, \dots, m} d(\underline{x}_i, C_k)$
  - Set  $b(i)=j$  \{  $b(i)$  is the index of the cluster that lies closet to  $\underline{x}_i$  \}
- End {for}
- For  $j=1$  to  $m$ 
  - Set  $C_j = \{\underline{x}_i \in X: b(i)=j\}$
  - If necessary, update representatives
- End {for}