# LECTURE 31

Unit IV – FP Growth Algorithm & Evaluation of
Association Patterns from Textbook

# LECTURE 32

Cluster Analysis
Source: Chapter 10, Data Mining:  Concepts and Techniques(3rd ed.)

# Cluster Analysis: Basic Concepts and Methods

- Cluster Analysis: Basic Concepts

- Partitioning Methods

- Hierarchical Methods

- Density-Based Methods

- ~~Grid-Based Methods~~

- Evaluation of Clustering

- Summary

# What is Cluster Analysis?

- Cluster: A collection of data objects
    - similar (or related) to one another within the same group
    - dissimilar (or unrelated) to the objects in other groups
- Cluster analysis (or *clustering*, *data segmentation, ...*)
    - Finding similarities between data according to the characteristics found in the data and grouping similar data objects into clusters
- Unsupervised learning: no predefined classes (i.e., *learning by observations* vs. learning by examples: supervised)
- Typical applications
    - As a stand-alone tool to get insight into data distribution
    - As a preprocessing step for other algorithms

# Applications of Cluster Analysis

- Data reduction
  - Summarization: Preprocessing for regression, PCA, classification, and association analysis
  - Compression: Image processing: vector quantization
- Hypothesis generation and testing
- Prediction based on groups
  - Cluster & find characteristics/patterns for each group
- Finding K-nearest Neighbors
  - Localizing search to one or a small number of clusters
- Outlier detection: Outliers are often viewed as those "far away" from any cluster

# Clustering: Application Examples

- Biology: taxonomy of living things: kingdom, phylum, class, order, family, genus and species

- Information retrieval: document clustering

- Land use: Identification of areas of similar land use in an earth observation database

- Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs

- City-planning: Identifying groups of houses according to their house type, value, and geographical location

- Earth-quake studies: Observed earthquake epicenters should be clustered along continent faults

- Climate: understanding earth climate, find patterns of atmospheric and ocean

- Economic Science: market resarch

# Basic Steps to Develop a Clustering Task

- Feature selection
  - Select info concerning the task of interest
  - Minimal information redundancy
- Proximity measure
  - Similarity of two feature vectors
- Clustering criterion
  - Expressed via a cost function or some rules
- Clustering algorithms
  - Choice of algorithms
- Validation of the results
  - Validation test (also, *clustering tendency* test)
- Interpretation of the results
  - Integration with applications

# Quality: What Is Good Clustering?

- A <u>good clustering</u> method will produce high quality clusters

    - high <u>intra-class</u> similarity: cohesive within clusters

    - low <u>inter-class</u> similarity: distinctive between clusters

- The <u>quality</u> of a clustering method depends on

    - the similarity measure used by the method

    - its implementation, and

    - Its ability to discover some or all the <u>hidden</u> patterns

# Measure the Quality of Clustering

- Dissimilarity/Similarity metric
    - Similarity is expressed in terms of a distance function, typically metric: $d(i, j)$
    - The definitions of distance functions are usually rather different for interval-scaled, boolean, categorical, ordinal ratio, and vector variables
    - Weights should be associated with different variables based on applications and data semantics
- Quality of clustering:
    - There is usually a separate "quality" function that measures the "goodness" of a cluster.
    - It is hard to define "similar enough" or "good enough"
        - The answer is typically highly subjective

# Considerations for Cluster Analysis

- Partitioning criteria

  - Single level vs. hierarchical partitioning (often, multi-level hierarchical partitioning is desirable)

- Separation of clusters

  - Exclusive (e.g., one customer belongs to only one region) vs. non-exclusive (e.g., one document may belong to more than one class)

- Similarity measure

  - Distance-based (e.g., Euclidian, road network, vector)  vs. connectivity-based (e.g., density or contiguity)

- Clustering space

  - Full space (often when low dimensional) vs. subspaces (often in high-dimensional clustering)

# Requirements and Challenges

- Scalability
  - Clustering all the data instead of only on samples
- Ability to deal with different types of attributes
  - Numerical, binary, categorical, ordinal, linked, and mixture of these
- Constraint-based clustering
  - User may give inputs on constraints
  - Use domain knowledge to determine input parameters
- Interpretability and usability
- Others
  - Discovery of clusters with arbitrary shape
  - Ability to deal with noisy data
  - Incremental clustering and insensitivity to input order
  - High dimensionality

# Major Clustering Approaches (I)

- Partitioning approach:
  - Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square errors
- Typical methods: k-means, k-medoids (Partitioning Around Medoids PAM) , CLARA (Clustering LARge Applications) CLARANS(Clustering Large Applications based upon RANdomized Search)

# Major Clustering Approaches (II)

- Hierarchical approach:

    - Create a hierarchical decomposition of the set of data (or objects) using some criterion

- Typical methods: DIANA (DIvisive ANAlysis), AGNES (AGglomerative NESting), BIRCH(Balanced Iterative Reducing and Clustering using Hierarchies), CAMELEON(Multiphase Hierarchical Clustering Using Dynamic Modeling)

- Density-based approach:

    - Based on connectivity and density functions

- Typical methods: DBSCAN(Density-Based Spatial Clustering of Applications with Noise) , OPTICS(Ordering Points to Identify the Clustering Structure), DenClue( Clustering Based on Density Distribution Functions

# Major Clustering Approaches (III)

- Grid-based approach:
    - based on a multiple-level granularity structure
    - Typical methods: STING, WaveCluster, CLIQUE

# Major Clustering Approaches (IV)

- Model-based:
  - A model is hypothesized for each of the clusters and tries to find the best fit of that model to each other
  - Typical methods: EM, SOM, COBWEB
- Frequent pattern-based:
  - Based on the analysis of frequent patterns
  - Typical methods: p-Cluster
- User-guided or constraint-based:
  - Clustering by considering user-specified or application-specific constraints
  - Typical methods: COD (obstacles), constrained clustering
- Link-based clustering:
  - Objects are often linked together in various ways
  - Massive links can be used to cluster objects: SimRank, LinkClus

# Cluster Analysis: Basic Concepts and Methods

- Cluster Analysis: Basic Concepts

- Partitioning Methods

- Hierarchical Methods

- Density-Based Methods

- ~~Grid-Based Methods~~

- Evaluation of Clustering

- Summary

# Partitioning Algorithms: Basic Concept

- Partitioning method: Partitioning a database **D** of **n** objects into a set of **k** clusters, such that the sum of squared distances is minimized (where $c_i$ is the centroid or medoid of cluster $C_i$)
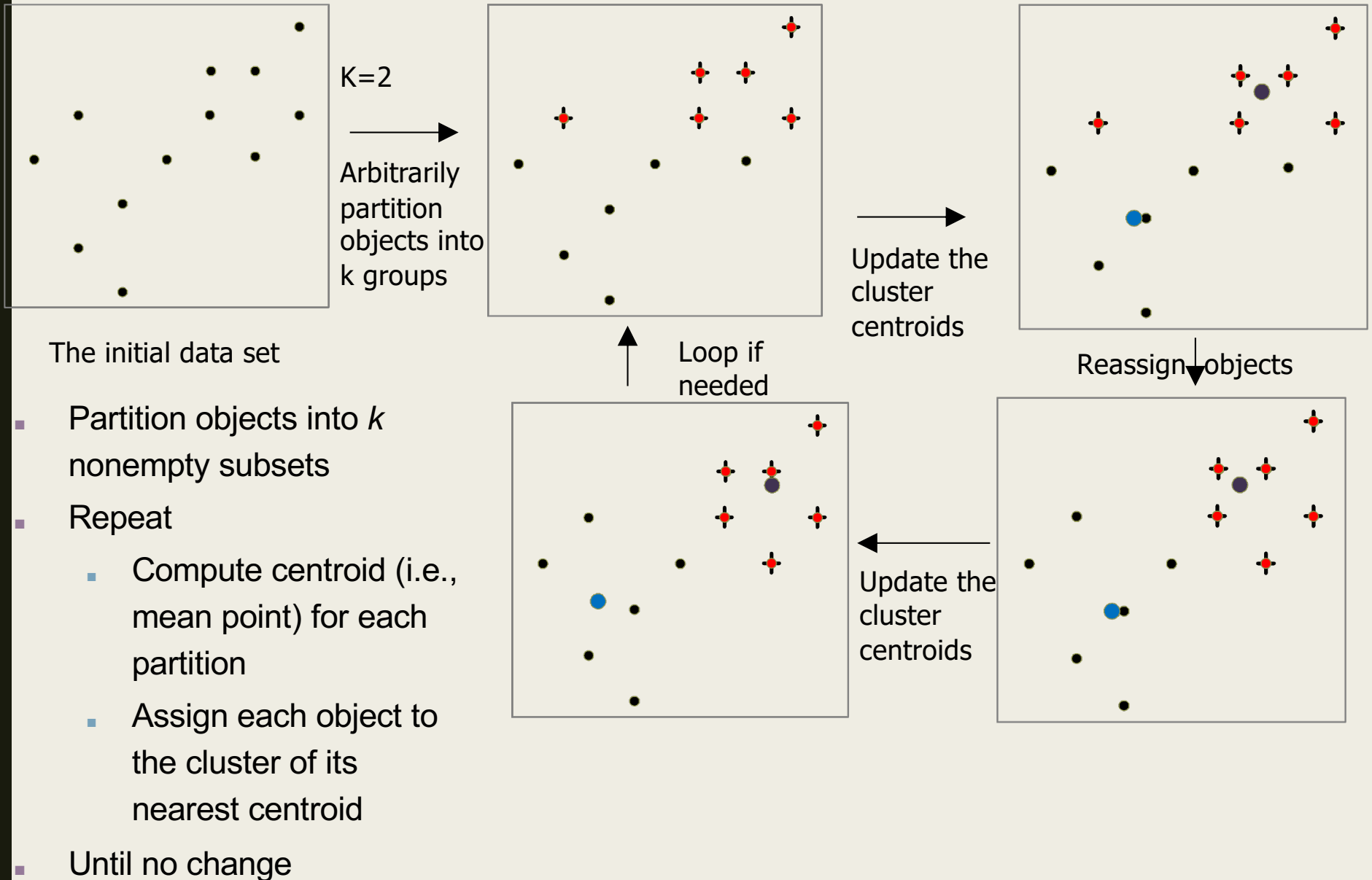
$$E = \Sigma_{i=1}^{k} \Sigma_{p \in C_i} (d(p, c_i))^2$$

- Given *k*, find a partition of *k clusters* that optimizes the chosen partitioning criterion

  - Global optimal: exhaustively enumerate all partitions

  - Heuristic methods: *k-means* and *k-medoids* algorithms

  - *k-means* (MacQueen'67, Lloyd'57/'82): Each cluster is represented by the center of the cluster

  - *k-medoids* or PAM (Partition around medoids) (Kaufman & Rousseeuw'87): Each cluster is represented by one of the objects in the cluster

# The *K-Means* Clustering Method

- Given *k*, the *k-means* algorithm is implemented in four steps:
    - Partition objects into *k* nonempty subsets
    - Compute seed points as the centroids of the clusters of the current partitioning (the centroid is the center, i.e., *mean point*, of the cluster)
    - Assign each object to the cluster with the nearest seed point
    - Go back to Step 2, stop when the assignment does not change

# An Example of *K-Means* Clustering

K=2

Arbitrarily partition objects into k groups

The initial data set

Update the cluster centroids

Reassign objects

Loop if needed

Update the cluster centroids

- Partition objects into *k* nonempty subsets
- Repeat
    - Compute centroid (i.e., mean point) for each partition
    - Assign each object to the cluster of its nearest centroid
- Until no change

48

# Comments on the *K-Means* Method

- Strength: *Efficient*: $O(tkn)$, where $n$ is # objects, $k$ is # clusters, and $t$ is # iterations. Normally, $k$, $t << n$.
    - Comparing: PAM: $O(k(n-k)^2)$, CLARA: $O(ks^2 + k(n-k))$
- Comment: Often terminates at a *local optimal*
- Weakness
    - Applicable only to objects in a continuous n-dimensional space
        - Using the k-modes method for categorical data
        - In comparison, k-medoids can be applied to a wide range of data
    - Need to specify $k$, the *number* of clusters, in advance (there are ways to automatically determine the best k (see Hastie et al., 2009)
    - Sensitive to noisy data and *outliers*
    - Not suitable to discover clusters with *non-convex shapes*

# Validity of clusters

- Why validity of clusters?

  - *Given some data, any clustering algorithm generates clusters*

  - *So, we need to make sure the clustering results are valid and meaningful.*

- Measuring the validity of clustering results usually involve

  - *Optimality of clusters*

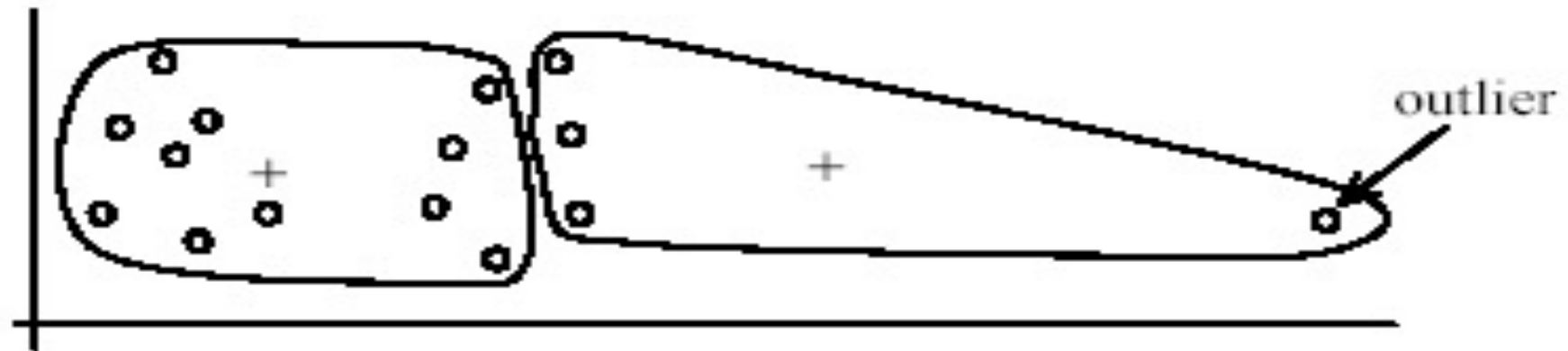  - *Verification of meaning of clusters*

# Optimality of clusters

- Optimal clusters should
  - *minimize distance **within** clusters (intracluster)*
  - *maximize distance **between** clusters (intercluster)*

- Example of intracluster measure
  - Squared error se

    *where $m_i$ is the mean of all instances in cluster $c_i$*

$$se = \sum_{i=1}^{k} \sum_{p \in c_i} \|p - m_i\|^2$$

# Weaknesses of k-means: Problems with outliers



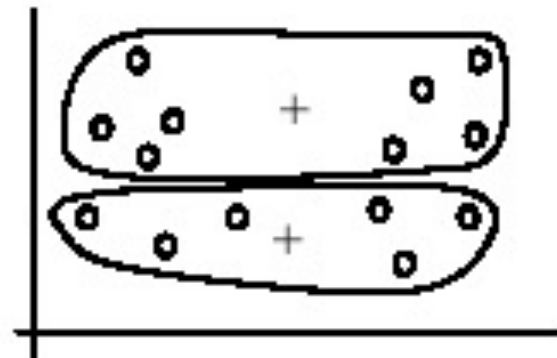(A): Undesirable clusters

(B): Ideal clusters

# Weaknesses of k-means (cont ...)

- ■ The algorithm is sensitive to initial seeds.
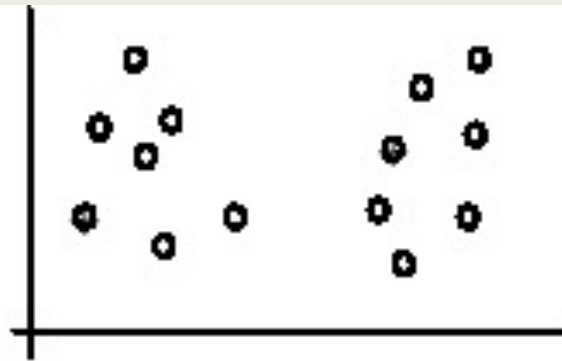


(A). Random selection of seeds (centroids)
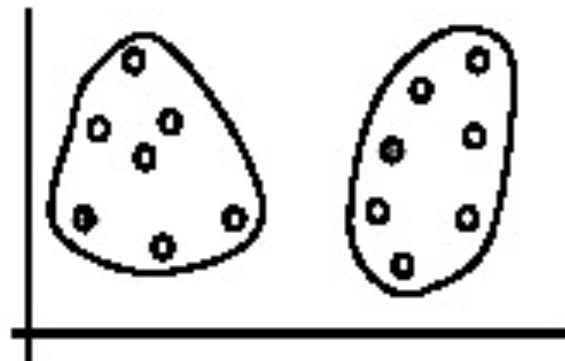
(B). Iteration 1
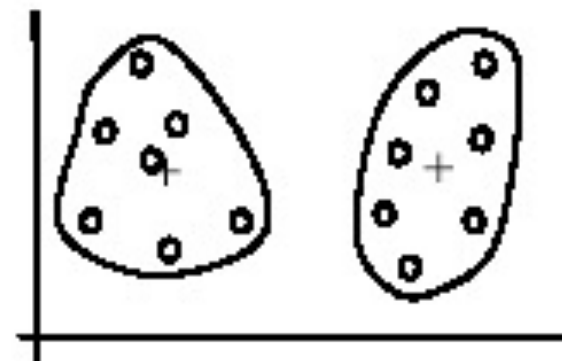
(C). Iteration 2

# Weaknesses of k-means (cont ...)

- If we use different seeds: good results



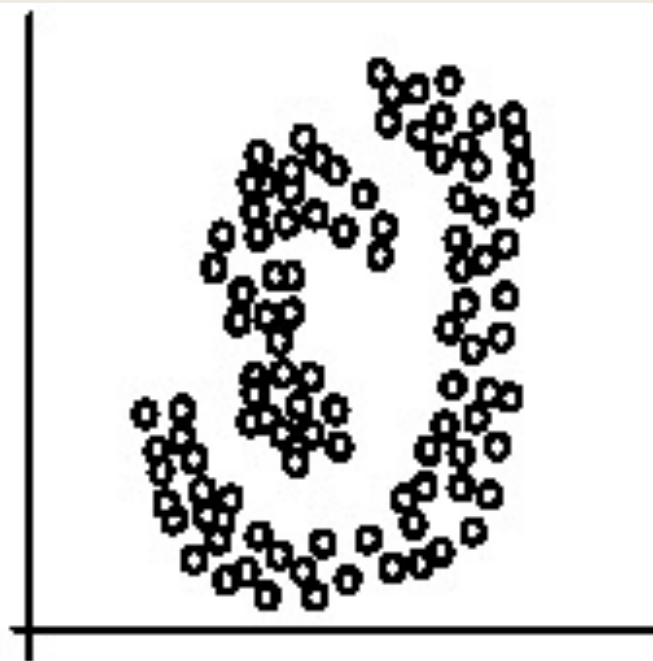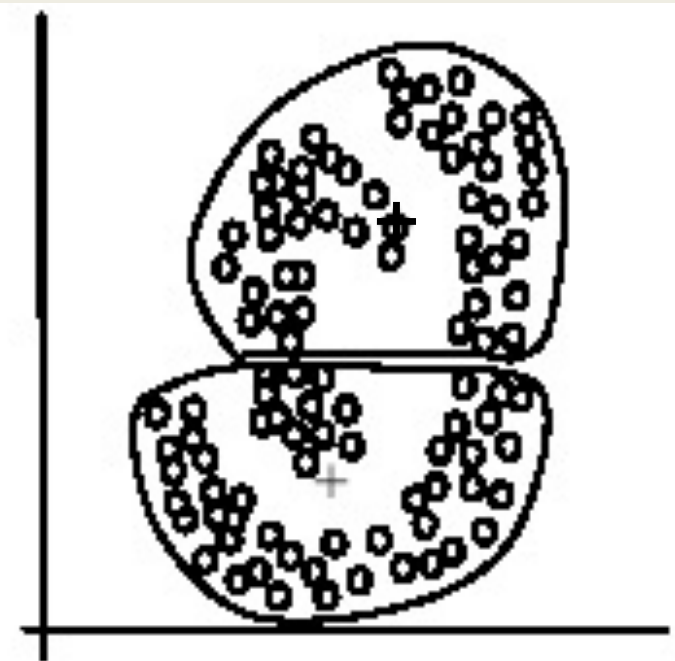(A). Random selection of $k$ seeds (centroids)

(B). Iteration 1

(C). Iteration 2

# Weaknesses of k-means (cont ...)

- The *k*-means algorithm is not suitable for discovering clusters that are not hyper-ellipsoids (or hyper-spheres).
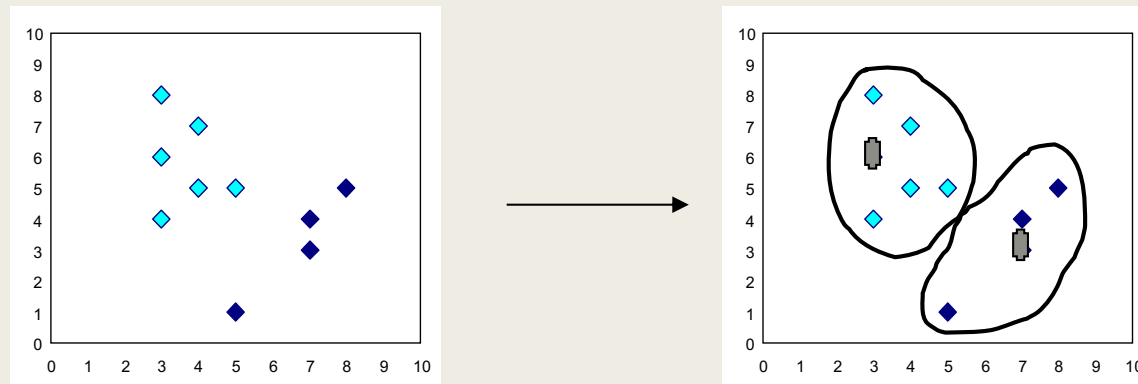


(A): Two natural clusters

(B): *k*-means clusters

# What is the problem of k-Means Method?

- The k-means algorithm is sensitive to outliers !

  - *Since an object with an extremely large value may substantially distort the distribution of the data.*

- K-Medoids:  Instead of taking the **mean** value of the object in a cluster as a reference point, **medoids** can be used, which is the **most centrally located** object in a cluster.

# Termination conditions

■ Several possibilities, e.g.,

– *A fixed number of iterations.*

– *Cluster partition unchanged.*
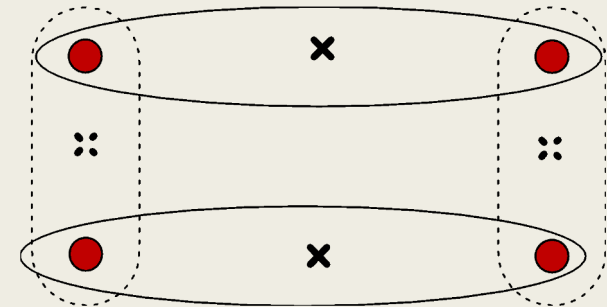
– *Centroid positions don't change.*

# *K*-means: summary

- Algorithmically, very simple to implement

- *K*-means converges, but it finds a local minimum of the cost function

- Works only for numerical observations

- *K* is a user input;

- Outliers can be considerable trouble to *K*-means

# Variations of the *K-Means* Method

- Most of the variants of the *k-means* which differ in

    - Selection of the initial *k* means

    - Dissimilarity calculations

    - Strategies to calculate cluster means

- Handling categorical data: *k-modes*

    - Replacing means of clusters with <u>modes</u>

    - Using new dissimilarity measures to deal with categorical objects

    - Using a <u>frequency</u>-based method to update modes of clusters

    - A mixture of categorical and numerical data: *k-prototype* method

# PAM: A Typical K-Medoids Algorithm

Total Cost = 20



Arbitrary choose k object as initial medoids

Assign each remaining object to nearest medoids

K=2

**Do loop**

**Until no change**

Total Cost = 26

Swapping O and O$_{ramdom}$

If quality is improved.

Compute total cost of swapping

Randomly select a nonmedoid object, O$_{ramdom}$