# Data Mining:
## Concepts and Techniques
### (3$^{rd}$ ed.)

## — Chapter 10 —

Jiawei Han, Micheline Kamber, and Jian Pei

University of Illinois at Urbana-Champaign &

Simon Fraser University

1

# Chapter 10. 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, 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 earth quake 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: <span style="color:red">cohesive</span> within clusters

  - low <u>inter-class</u> similarity: <span style="color:red">distinctive</span> 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 of 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

# Requirements and Challenges

- 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
  - To choose the locations for a given number of new automatic teller machines (ATMs) in a city.
- Interpretability and usability
  - Users want clustering results to be interpretable, comprehensible, and usable.
  - clustering may need to be tied in with specific semantic interpretations and applications.
  - It is important to study how an application goal may influence the selection of clustering features and clustering methods.

# Requirements and Challenges

- Scalability
    - Clustering all the data instead of only on samples
- Discovery of clusters with arbitrary shape
    - Many clustering algorithms determine clusters based on Euclidean or Manhattan distance measures (Chapter 2).
    - Algorithms based on such distance measures tend to find spherical clusters with similar size and density.
    - It is important to develop algorithms that can detect clusters of arbitrary shape.
- Ability to deal with noisy data
    - Most real-world data sets contain outliers and/or missing, unknown, or erroneous data.
        - For example Sensor readings are often noisy—some readings may be inaccurate due to the sensing mechanisms.

# Requirements and Challenges

- Incremental clustering and insensitivity to input order
  - Some clustering algorithms cannot incorporate incremental updates into existing clustering structures and, instead, have to recompute a new clustering from scratch.
  - Clustering algorithms may also be sensitive to the input data order.
  - Clustering algorithms may return dramatically different clusterings depending on the order in which the objects are presented.
  - Incremental clustering algorithms and algorithms that are insensitive to the input order are needed.
- High dimensionality
  - A data set can contain numerous dimensions or attributes.
  - Clustering algorithms are good at handling low-dimensional data such as data sets involving only two or three dimensions.
  - Finding clusters of data objects in a highdimensional space is challenging, especially considering that such data can be very sparse and highly skewed.

# 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)

# Considerations for Cluster Analysis

- Similarity measure

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

- Clustering space

  - Many clustering methods search for clusters within the entire given data space.

  - These methods are useful for low-dimensionality data sets.

  - With highdimensional data, however, there can be many irrelevant attributes, which can make similarity measurements unreliable.

  - It's better to instead search for clusters within different subspaces of the same data set.

# 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, CLARANS
- Hierarchical approach:
  - Create a hierarchical decomposition of the set of data (or objects) using some criterion
  - Typical methods: Diana, Agnes, BIRCH, CAMELEON
- Density-based approach:
  - Based on connectivity and density functions
  - Typical methods: DBSACN, OPTICS, DenClue
- Grid-based approach:
  - based on a multiple-level granularity structure
  - Typical methods: STING, WaveCluster, CLIQUE

| Method | General Characteristics |
|---|---|
| Partitioning methods | – Find mutually exclusive clusters of spherical shape<br>– Distance-based<br>– May use mean or medoid (etc.) to represent cluster center<br>– Effective for small- to medium-size data sets |
| Hierarchical methods | – Clustering is a hierarchical decomposition (i.e., multiple levels)<br>– Cannot correct erroneous merges or splits<br>– May incorporate other techniques like microclustering or consider object "linkages" |
| Density-based methods | – Can find arbitrarily shaped clusters<br>– Clusters are dense regions of objects in space that are separated by low-density regions<br>– Cluster density: Each point must have a minimum number of points within its "neighborhood"<br>– May filter out outliers |
| Grid-based methods | – Use a multiresolution grid data structure<br>– Fast processing time (typically independent of the number of data objects, yet dependent on grid size) |

# Chapter 10. 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

- <u>Partitioning method</u>: 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 $C_i$)
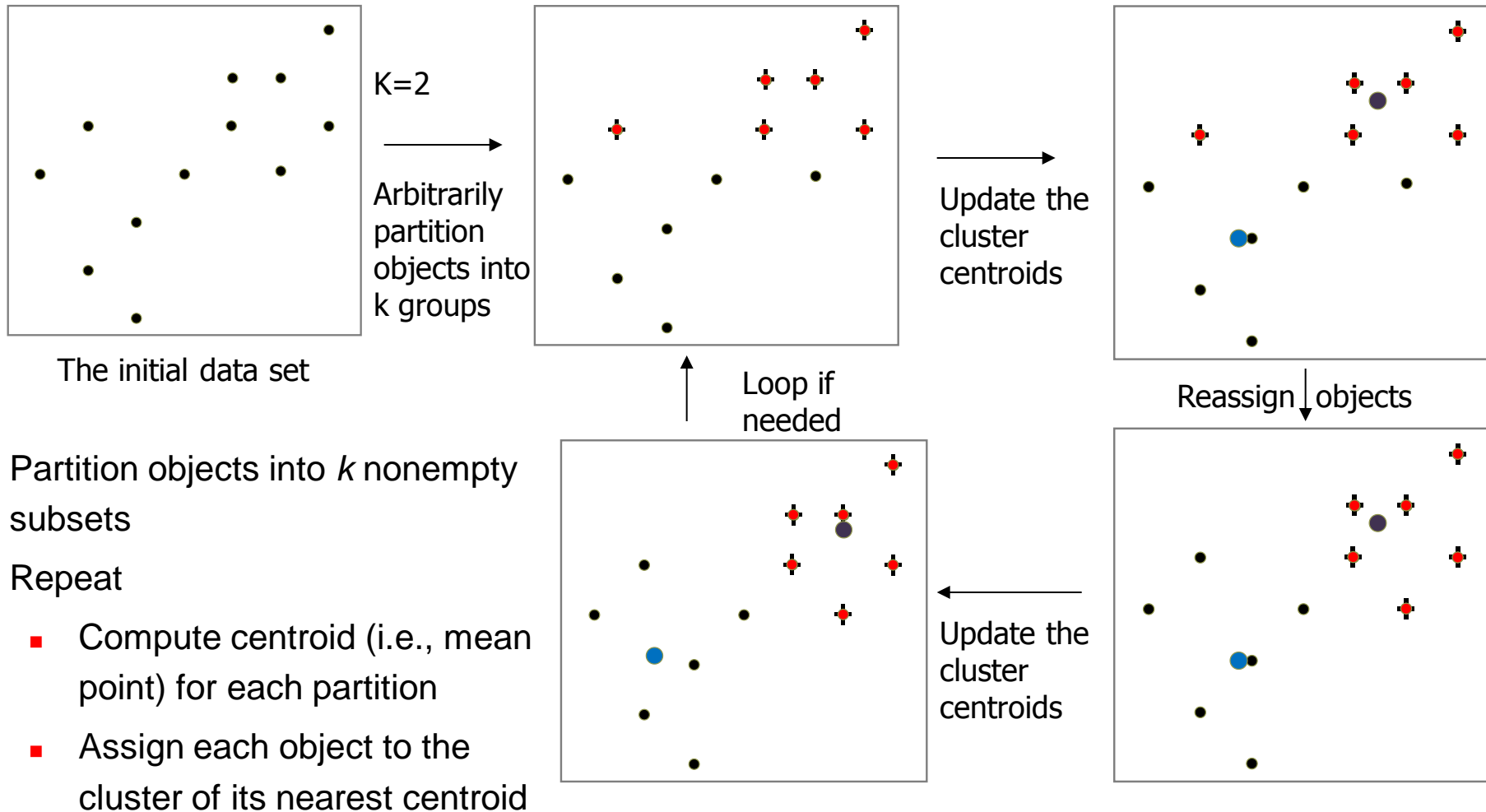
$$E = \Sigma_{i=1}^{k} \Sigma_{p \in C_i} (d(p, c_i))^2$$

- 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

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

  - <u>*k-medoids*</u> 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
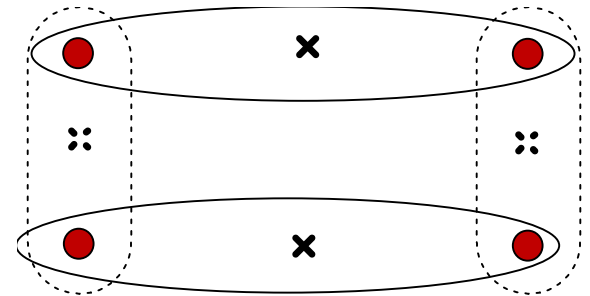
# An Example of *K-Means* Clustering

The initial data set

K=2

Arbitrarily partition objects into k groups

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

# 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*
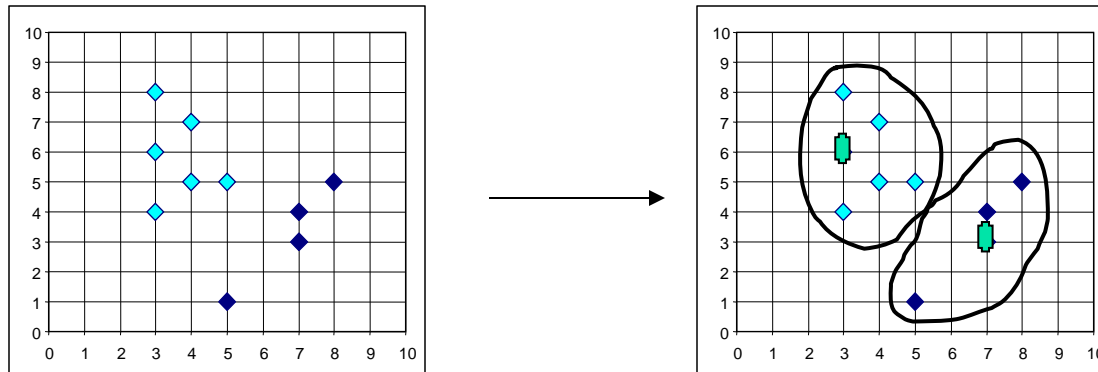    - Not suitable to discover clusters with *non-convex shapes*

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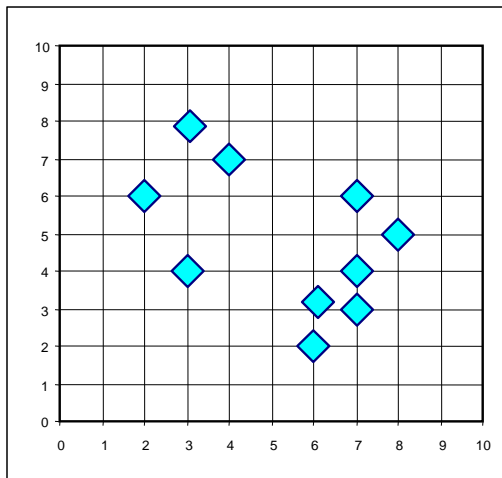
  - A mixture of categorical and numerical data: *k-prototype* method

# What Is the Problem of the 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
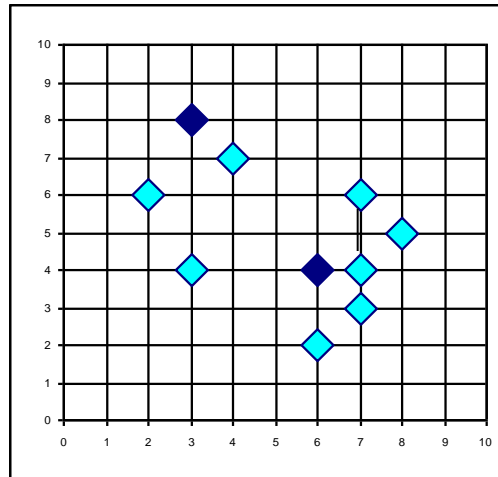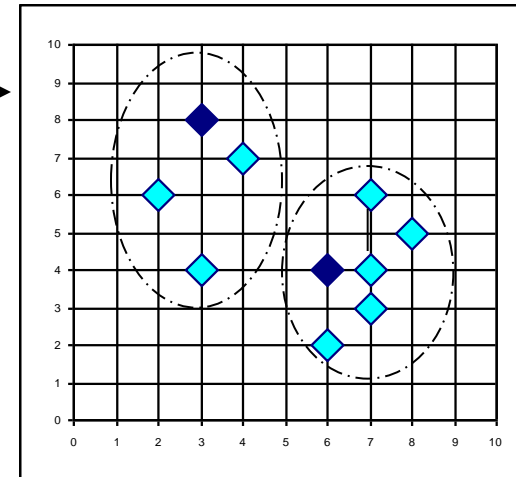
# PAM: A Typical K-Medoids Algorithm

Total Cost = 20
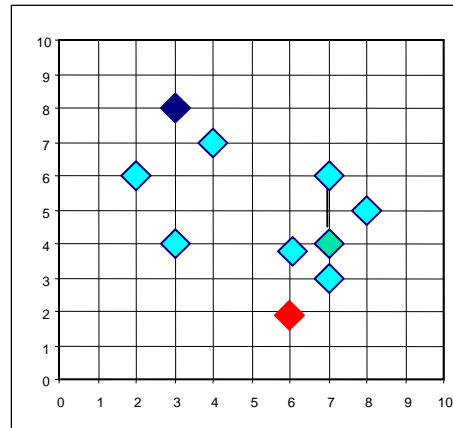


K=2

Arbitrary choose k object as initial medoids
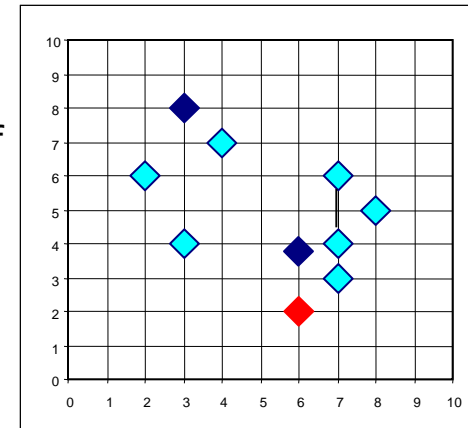
Assign each remaining object to nearest medoids

Randomly select a nonmedoid object,$O_{ramdom}$

**Do loop**

**Until no change**

Swapping O and $O_{ramdom}$

If quality is improved.

Total Cost = 26

Compute total cost of swapping

# The K-Medoid Clustering Method

- *K-Medoids* Clustering: Find *representative* objects (<u>medoids</u>) in clusters

  - *PAM* (Partitioning Around Medoids, Kaufmann & Rousseeuw 1987)

    - Starts from an initial set of medoids and iteratively replaces one of the medoids by one of the non-medoids if it improves the total distance of the resulting clustering

    - *PAM* works effectively for small data sets, but does not scale well for large data sets (due to the computational complexity)

- Efficiency improvement on PAM

  - *CLARA* (Kaufmann & Rousseeuw, 1990): PAM on samples

  - *CLARANS* (Ng & Han, 1994): Randomized re-sampling