

## 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

19

## What is Cluster Analysis?

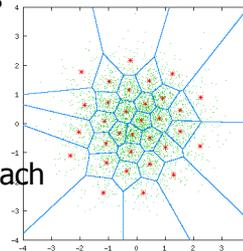
---

- 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

20

## 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



21

## 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 research

22

## 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

23

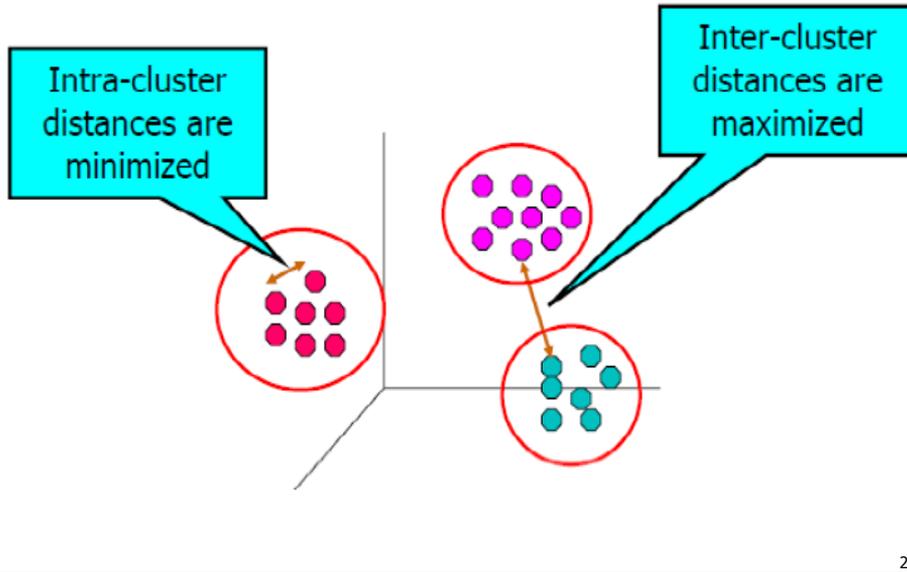
## Quality: What Is Good Clustering?

---

- A good clustering method will produce high quality clusters
  - high intra-class similarity: **cohesive** within clusters
  - low inter-class similarity: **distinctive** between clusters
- The quality 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 hidden patterns

24

## Good clustering



## 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)

27

## 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

28

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

29

## Major Clustering Approaches (II)

---

- 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

30

## 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

31

## 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  $C_i$ )

$$E = \sum_{i=1}^k \sum_{p \in C_i} (d(p, c_i))^2$$

- 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

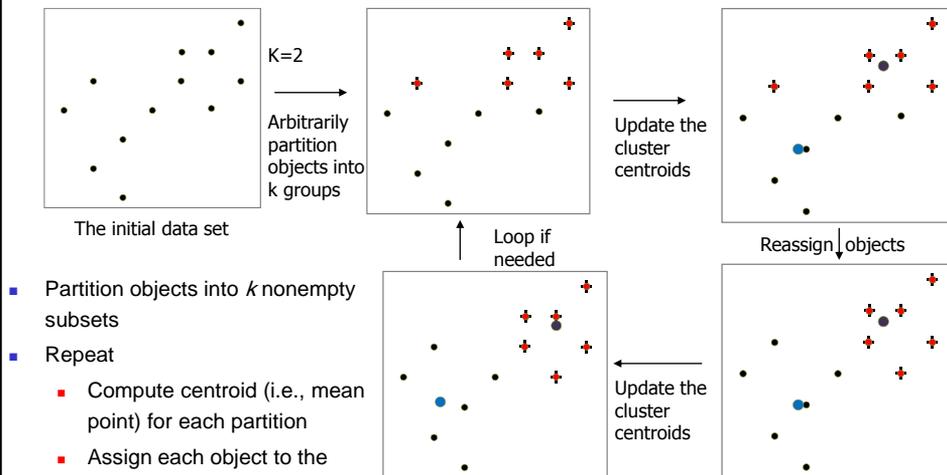
32

# 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

33

# An Example of *K-Means* Clustering



- 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

34

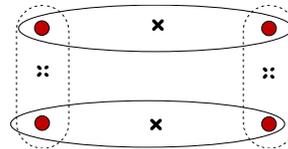
## Comments on the *K-Means* Method

- Strength: *Efficient*:  $O(tkn)$ , where  $n$  is # objects,  $k$  is # clusters, and  $t$  is # iterations. Normally,  $k, t \ll n$ .
  - 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*

35

## Variations of the *K-Means* Method

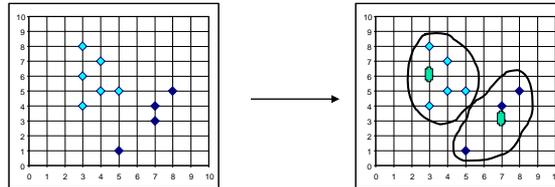
- Most of the variants of the *k-means* which differ in
  - Selection of the initial  $k$  means
  - Dissimilarity calculations
  - Strategies to calculate cluster means
- Handling categorical data: *k-modes*
  - Replacing means of clusters with modes
- Using new dissimilarity measures to deal with categorical objects
- Using a frequency-based method to update modes of clusters
- A mixture of categorical and numerical data: *k-prototype* method



36

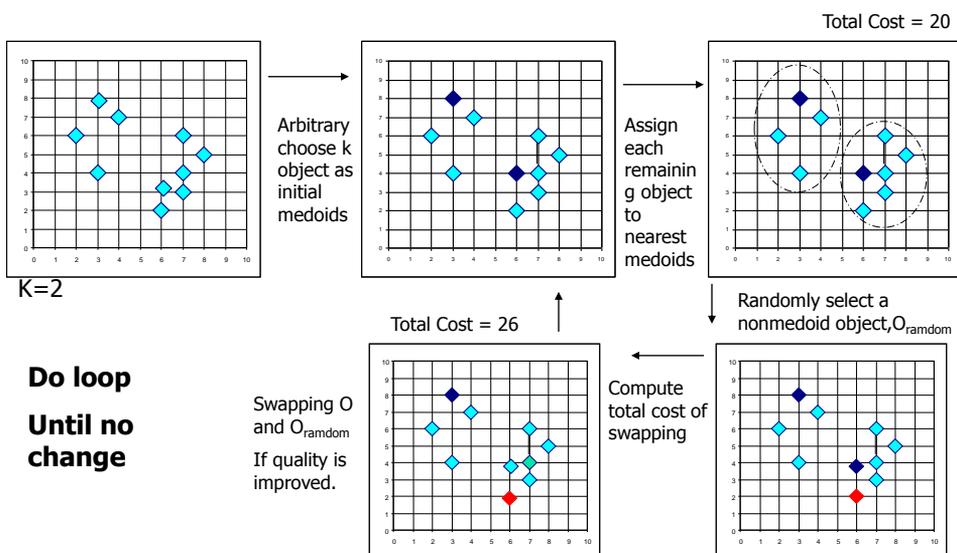
## What Is the Problem of the K-Means Method?

- 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



37

## PAM: A Typical K-Medoids Algorithm



## The K-Medoid Clustering Method

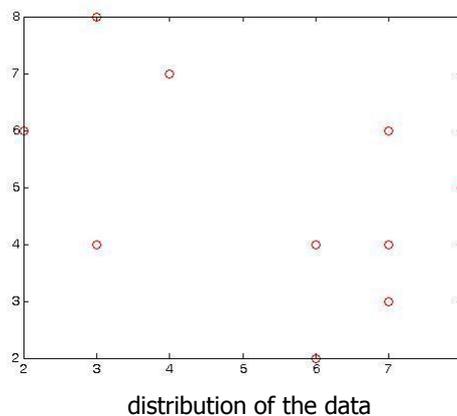
- *K-Medoids* Clustering: Find *representative* objects (medoids) in clusters
  - *PAM* (Partitioning Around Medoids, Kaufmann & Rousseeuw 1987)
    - Starts from an initial set of medoids and iteratively replaces one of the medoids by one of the non-medoids if it improves the total distance of the resulting 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 (**Projects for students**)
    - *CLARA* (Kaufmann & Rousseeuw, 1990): PAM on samples
    - *CLARANS* (Ng & Han, 1994): Randomized re-sampling

39

## Example: k-medoids

- Cluster the following data set of ten objects into two clusters i.e.  $k = 2$ .

	F <sub>1</sub>	F <sub>2</sub>
X <sub>1</sub>	2	6
X <sub>2</sub>	3	4
X <sub>3</sub>	3	8
X <sub>4</sub>	4	7
X <sub>5</sub>	6	2
X <sub>6</sub>	6	4
X <sub>7</sub>	7	3
X <sub>8</sub>	7	4
X <sub>9</sub>	8	5
X <sub>10</sub>	7	6



40

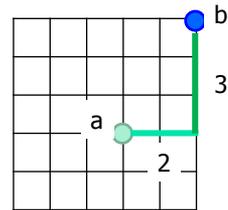
## Example: k-medoids

■ **1<sup>st</sup> step:** Initialize  $k$  centers.

- Let us assume  $x_2$  and  $x_8$  are selected as medoids, so the centers are  $c_1 = (3,4)$  and  $c_2 = (7,4)$
- Calculate distances to each center so as to associate each data object to its nearest medoid. Cost is calculated using **Manhattan distance (Minkowski distance metric with  $r = 1$ )**.
- **Manhattan distance ( $L_1$ -distance)**

$$d_1(\mathbf{p}, \mathbf{q}) = \|\mathbf{p} - \mathbf{q}\|_1 = \sum_{i=1}^n |p_i - q_i|,$$

$$d_m(x,y) = |x_a - x_b| + |y_a - y_b|$$



41

## Example: k-medoids

Costs to the nearest medoid are shown **bold** in the table

Cost (distance) to $c_1$					Cost (distance) to $c_2$				
$i$	$c_1$		Data objects ( $X_i$ )	Cost (distance)	$i$	$c_2$		Data objects ( $X_i$ )	Cost (distance)
$X_1$	3	4	2 6	<b>3</b>	$X_1$	7	4	2 6	7
$X_3$	3	4	3 8	<b>4</b>	$X_3$	7	4	3 8	8
$X_4$	3	4	4 7	<b>4</b>	$X_4$	7	4	4 7	6
$X_5$	3	4	6 2	5	$X_5$	7	4	6 2	<b>3</b>
$X_6$	3	4	6 4	3	$X_6$	7	4	6 4	<b>1</b>
$X_7$	3	4	7 3	5	$X_7$	7	4	7 3	<b>1</b>
$X_9$	3	4	8 5	6	$X_9$	7	4	8 5	<b>2</b>
$X_{10}$	3	4	7 6	6	$X_{10}$	7	4	7 6	<b>2</b>

cost between any two points is found using formula

$$\text{cost}(x, c) = \sum_{i=1}^d |x_i - c_i|$$

42

## Example: k-medoids

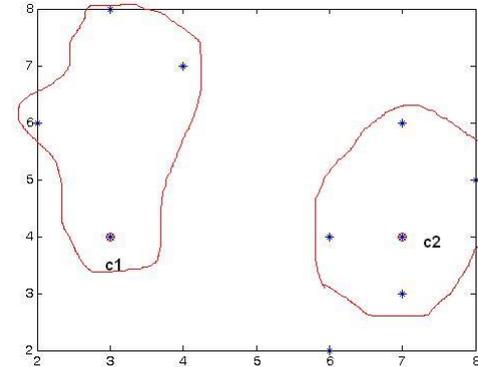
Since the points (2,6) (3,8) and (4,7) are closer to  $c_1$  hence they form one cluster whilst remaining points form another cluster.

Then the clusters become:

$$\text{Cluster}_1 = \{(3,4)(2,6)(3,8)(4,7)\}$$

$$\text{Cluster}_2 = \{(7,4)(6,2)(6,4)(7,3)(8,5)(7,6)\}$$

clusters after step 1

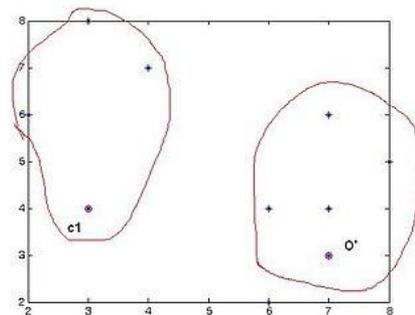


42

## Example: k-medoids

■ **2<sup>nd</sup> step:** Select one of the non-medoids  $O'$

- Let us assume  $O' = (7,3)$ , i.e.  $x_7$ .
- So now the medoids are  $c_1(3,4)$  and  $O'(7,3)$
- If  $c_1$  and  $O'$  are new medoids, calculate the total cost involved



clusters after step 2

44

## Example: k-medoids

- By using the formula in the step 1

$i$	$c_1$	Data objects ( $X_i$ )	Cost (distance)
1	3 4	2 6	3
3	3 4	3 8	4
4	3 4	4 7	4
5	3 4	6 2	5
6	3 4	6 4	3
8	3 4	7 4	4
9	3 4	8 5	6
10	3 4	7 6	6

$i$	$O'$	Data objects ( $X_i$ )	Cost (distance)
1	7 3	2 6	8
3	7 3	3 8	9
4	7 3	4 7	7
5	7 3	6 2	2
6	7 3	6 4	2
8	7 3	7 4	1
9	7 3	8 5	3
10	7 3	7 6	3

- Total cost =  $3+4+4+2+2+1+3+3=22$
- So cost of swapping medoid from  $c_2$  to  $O'$  is  
 $S = \text{current cost} - \text{past cost}$   
 $= 22 - 20$   
 $= 2 > 0$

45

## Example: k-medoids

- So cost of swapping medoid from  $c_2$  to  $O'$  is  
 $S = \text{current cost} - \text{past cost}$   
 $= 22 - 20$   
 $= 2 > 0$
- So moving to  $O'$  would be a bad idea, so the previous choice was good. So we try other non-medoids and found that our first choice was the best. So the configuration does not change and algorithm terminates here (i.e. there is no change in the medoids).
- It may happen some data points may shift from one cluster to another cluster depending upon their closeness to medoid.

46