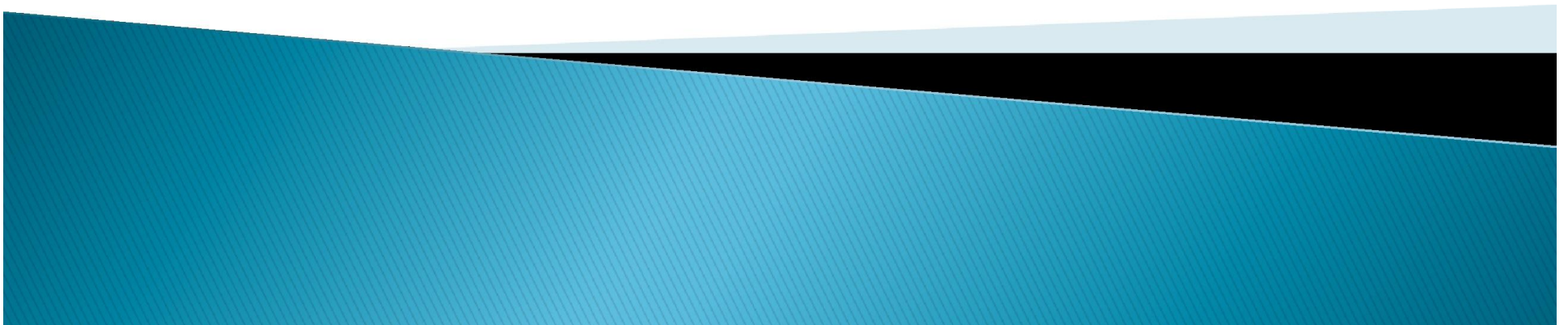# Cluster Analysis (Part A)

## Nazerfard, Ehsan

nazerfard@aut.ac.ir

# Cluster Analysis

❑ Cluster analysis or clustering: Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups.
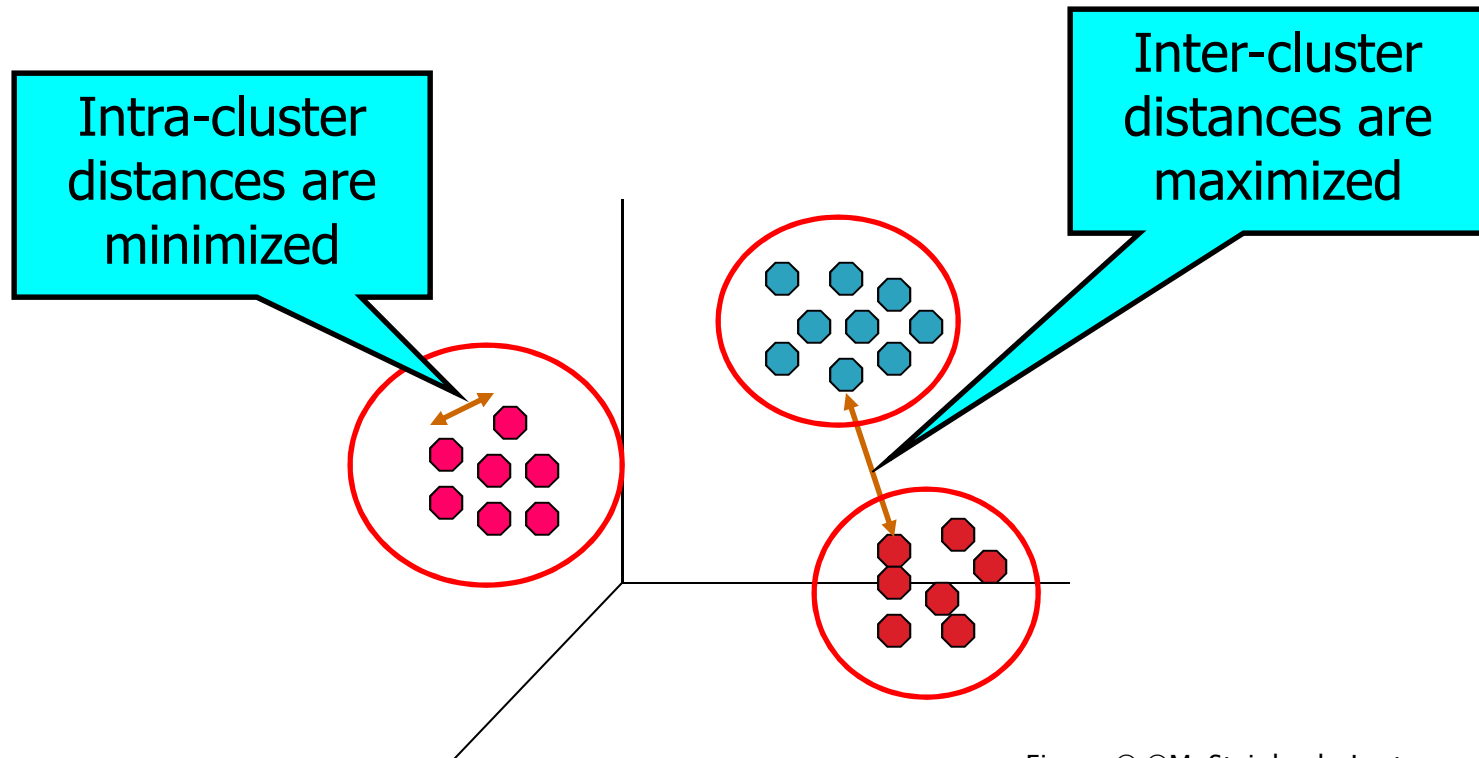
Intra-cluster distances are minimized

Inter-cluster distances are maximized

Figure © ©M. Steinbach, Lecture on Clustering

# Applications

- Document clustering (news, …)
- Sentiment analysis (customer reviews, …)
- Gene expression clustering
- Clustering of patients based on phenotypic and genotypic factors for efficient disease diagnosis
- Market Segmentation
- Anomaly detection
- Fraud detection
- Finding groups of driver behaviors based upon patterns of automobile motions (normal, drunken, sleepy, rush hour driving, etc.)
- …

# Major Clustering Approaches

- **Partitioning-based 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, CLARA, CLARANS
- **Density-based approach**
  - Based on connectivity and density functions
  - Typical methods: DBSACN, OPTICS, DenClue
- **Hierarchical approach**
  - Create a hierarchical decomposition of the set of data (or objects) using some criterion
  - Typical methods: Agnes, Diana, BIRCH, CURE, CAMELEON
- **Model-based approach**
  - 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

# Major Clustering Approaches (cont.)

- ❑ Grid-based approach
  - o Based on a multiple-level granularity structure
  - o Typical methods: STING, CLIQUE, WaveCluster
- ❑ Frequent Pattern-based approach
  - o Based on the analysis of frequent patterns
  - o Typical methods: p-Cluster
- ❑ Support Vector approach
  - o Based on the idea of mapping data points into higher dimensional feature space via a kernel function.
  - o Typical methods: SVC, Kernel K-means
- ❑ Graph Theoretic approach
  - o Typical methods: Spectral Clustering
- ❑ …

# Partitioning-based Approach

❑ Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square errors.
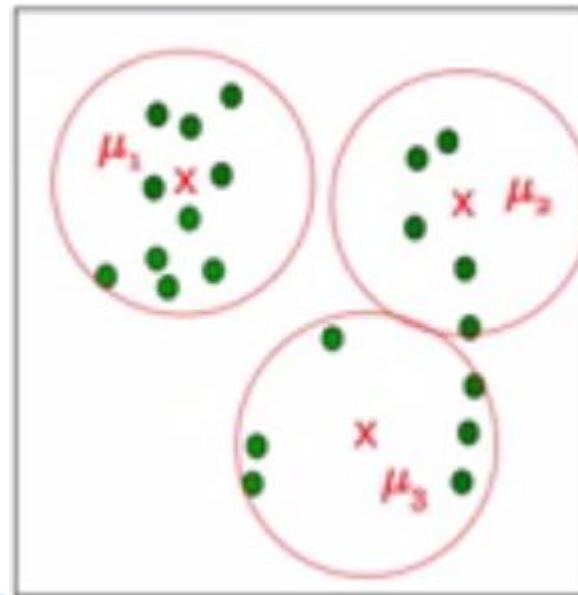
➢ Example: K-means

# K-means Clustering

❑ Assume $K$ clusters

❑ Iterate between two following steps:
  ○ Updating the assignment of data to clusters
  ○ Updating the cluster's summarization
  ___
  ○ sklearn.cluster.KMeans
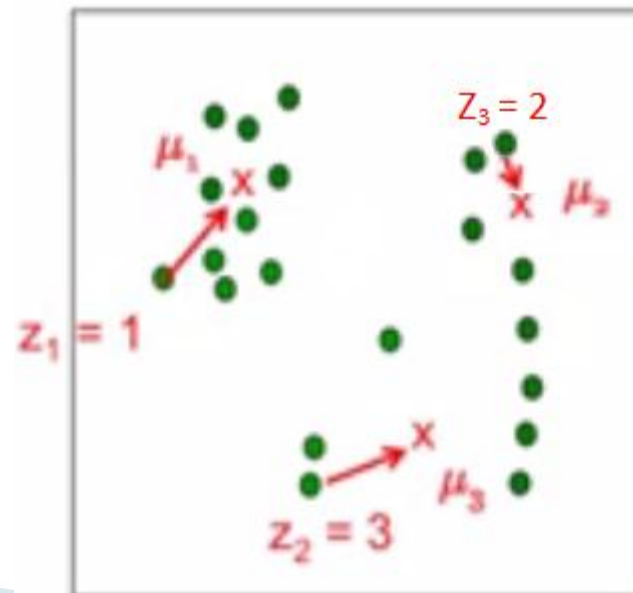
# K-means Clustering (cont.)

- ❏ Assume $K$ clusters
- ❏ Iterate between two following steps
  - A. Updating the assignment of data to clusters
  - B. Updating the cluster's summarization
- ❏ Each cluster $C$ is described by a centroid $\boldsymbol{\mu_c}$

# K-means Clustering (cont.)

- Assume $K$ clusters
- Iterate between two following steps:
  - A. Updating the assignment of data to clusters
  - B. Updating the cluster's summarization
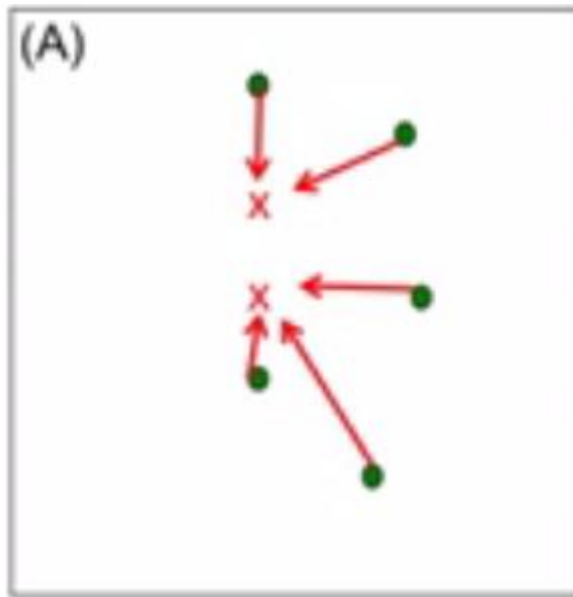- Assignment of $i_{th}$ example: $z_i \in 1..K$

# K-means Clustering (cont.)

❑ Iterate until convergence

   A. For each data, find the closest centroid:

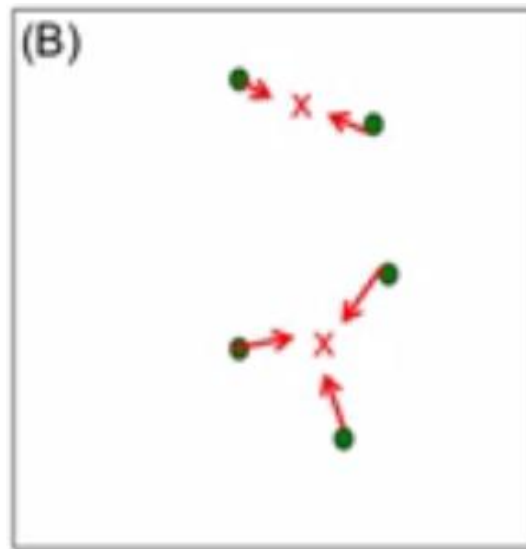$$z_i = \underset{c}{\mathrm{argmin}} ||\boldsymbol{x_i} - \boldsymbol{\mu_c}||^2 , \forall i$$



(A)

# K-means Clustering (cont.)

❑ Iterate until convergence
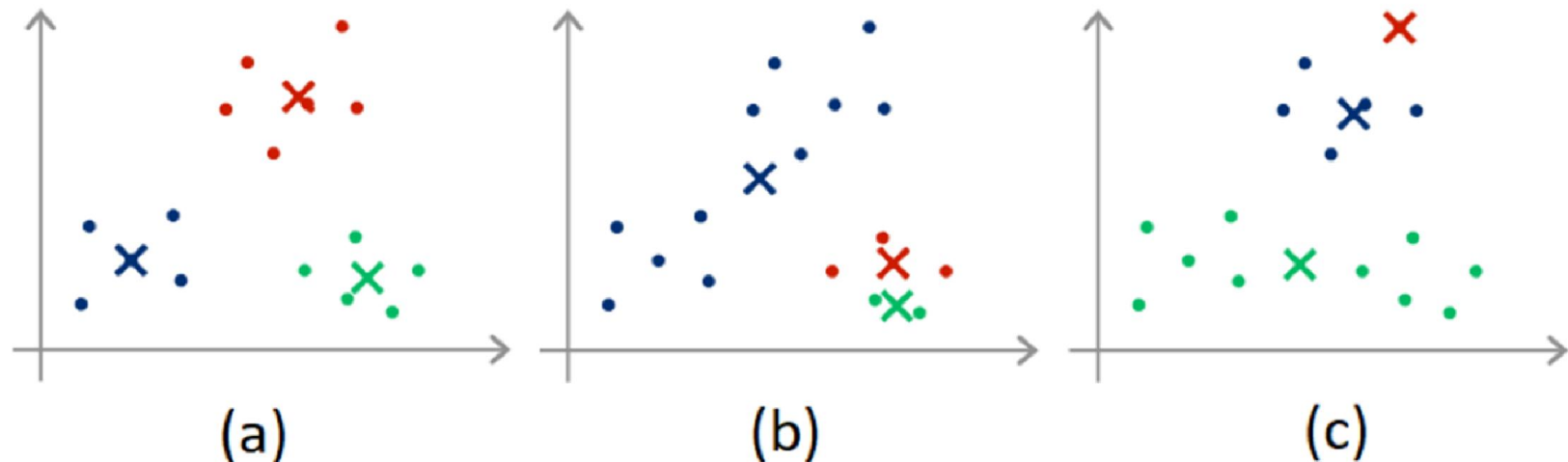
B. Set each cluster to the mean of all assigned data:

$$\forall c, \qquad \boldsymbol{\mu_c} = 1/m_c \sum_{i \in S_c} x_i \qquad\qquad S_c = \{i: z_i = c\}, m_c = |S_c|$$



(B)

Demo

# K-means Properties

❑ Poor initialization may lead to poor clustering



(a)       (b)       (c)

❑ Solution?
  o Multiple Initializations → randomness
  o K-means++, Intelligent K-means

# K-means Properties (cont.)

❑ Distance metrics
- $l_1$ norm (Manhattan distance)
- $l_2$ norm (Euclidean distance)
- Cosine similarity

❑ Centroids
- Mean
- Median → Outliers?
- Medoid
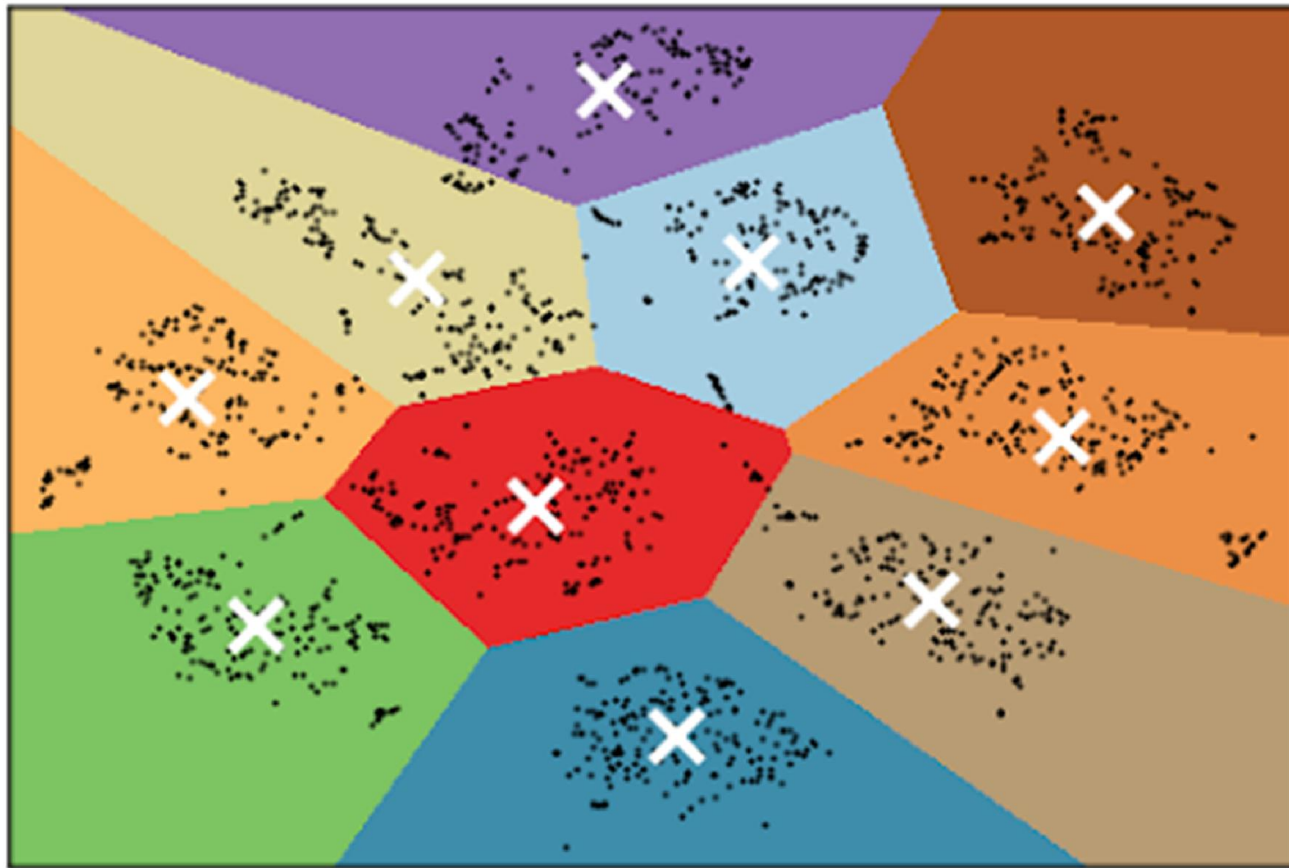  - Most commonly used on data when a mean or centroid cannot be defined, such as graphs.
- …

# K-means Properties (cont.)

- Instance-based

- Time complexity: $O(tkm)$

- Non-parametric

- Linearly separable data

# K-means: Linear Separable

# Sum of Square Error

❑ Sum of Square Error (SSE)
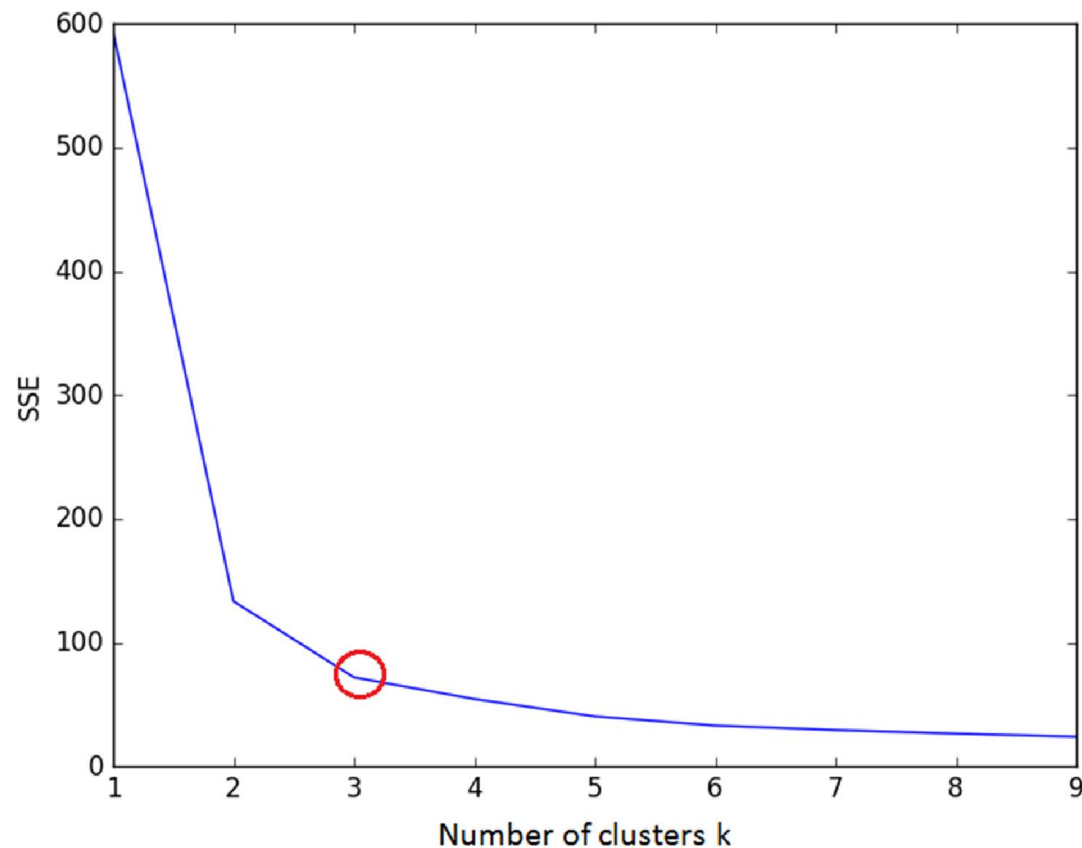
$$SSE = \sum_k \sum_{x_i \in C_k} ||x_i - C_k||^2$$

- Goal: minimizing within-cluster distance
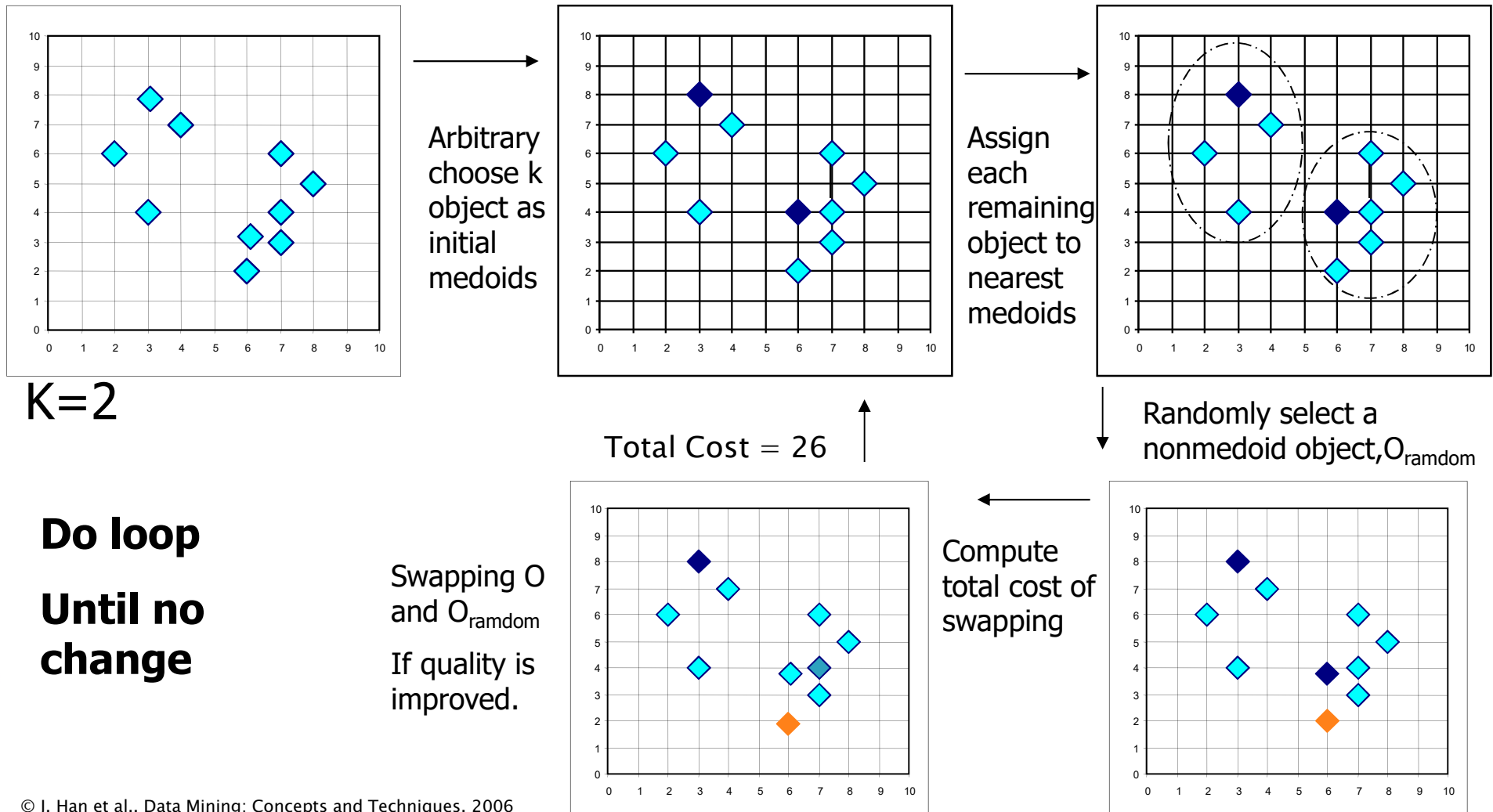
# Optimal number of Clusters

❑ Elbow method

# K-means Variations

- ❑ K-medoids/PAM (Partitioning Around Medoids)

- ❑ CLARA (Clustering Large Applications)

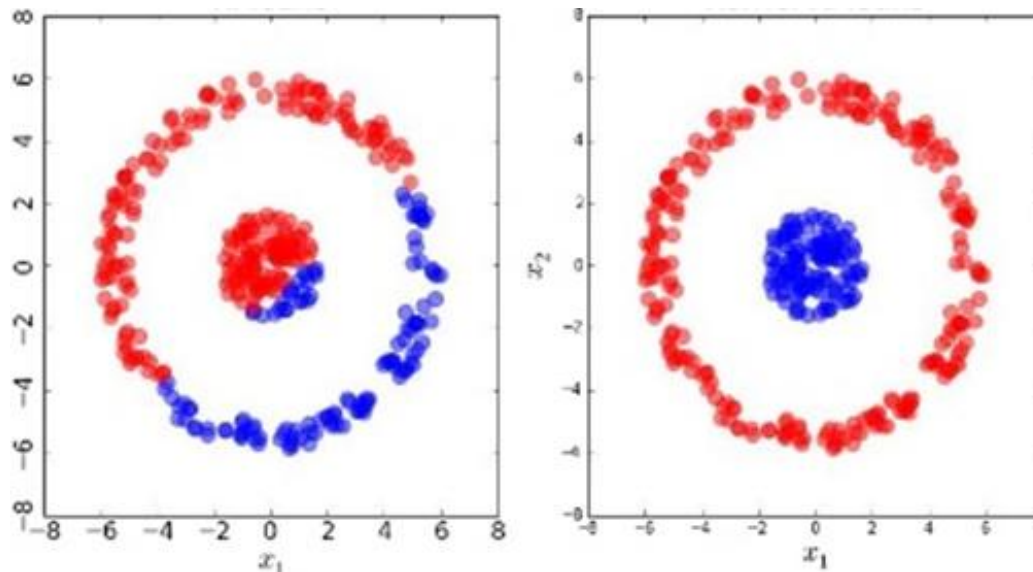- ❑ CLARANS (A Clustering Algorithm based on Randomized Search)

# PAM Algorithm



Total Cost = 20

K=2

Arbitrary choose k object as initial medoids

Assign each remaining object to nearest medoids

**Do loop**

**Until no change**

Swapping O and O$_{ramdom}$

If quality is improved.

Total Cost = 26

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

Compute total cost of swapping

© J. Han et al., Data Mining: Concepts and Techniques, 2006

# K-means Variations (cont.)

❑ Fuzzy C-means

❑ Kernel K-means

# Further Reading

- Mean Shift Clustering

- Clustering Categorical Data
  - ROCK (robust clustering algorithm for categorical attributes)
  - Sudipto Guha, Rajeev Rastogi, Kyuseok Shim, ICDE'99

# References

- Jiawei Han, Micheline Kamber and Jian Pei, Data Mining: Concepts and Techniques, 3$^{rd}$ edition, 2006.