# Kmeans Clustering

## Nazerfard, Ehsan

nazerfard@aut.ac.ir
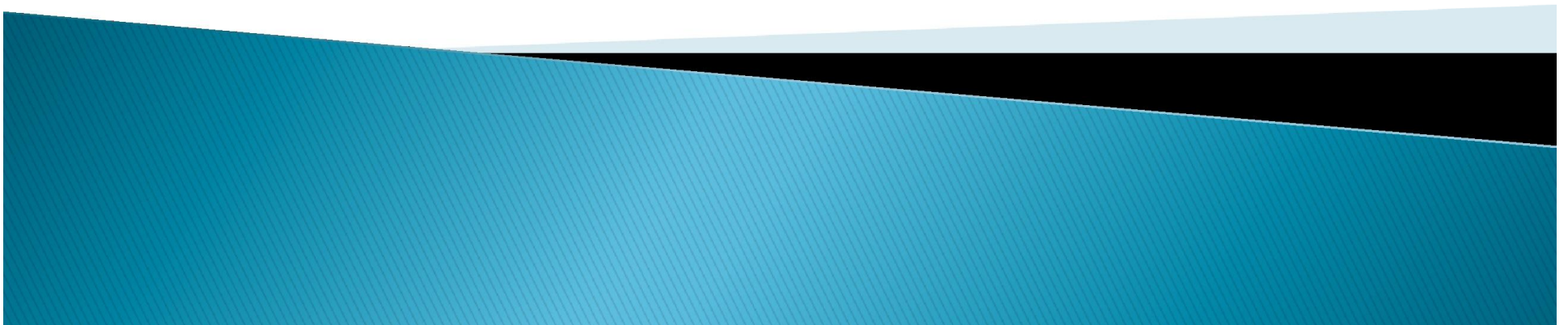
# Cluster Analysis

❑ 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.
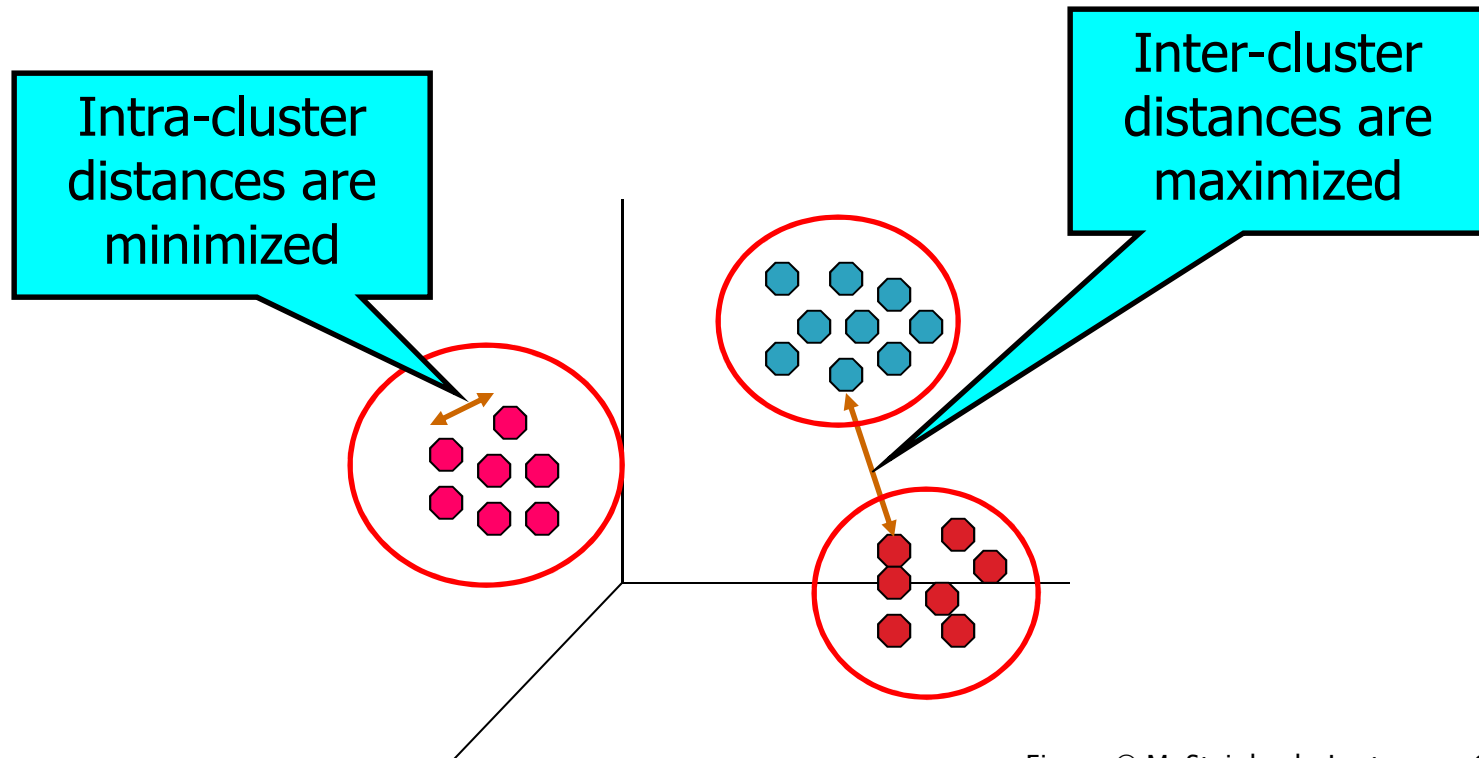
Intra-cluster distances are minimized

Inter-cluster distances are maximized

Figure © M. Steinbach, Lecture on Clustering

# Applications

- Document clustering (news, …)
- Community Detection in Social Networks
- 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**
  - o Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square errors.
  - o Typical methods: <span style="color:red">k-means</span>, k-medoids, CLARA, CLARANS
- ❑ **Density-based approach**
  - o Based on connectivity and density functions
  - o Typical methods: <span style="color:red">DBSCAN</span>, OPTICS, DenClue
- ❑ **Hierarchical approach**
  - o Create a hierarchical decomposition of the set of data (or objects) using some criterion
  - o Typical methods: <span style="color:red">Agnes, Diana</span>, BIRCH, CURE, CHAMELEON
- ❑ **Model-based approach**
  - o A model is hypothesized for each of the clusters and tries to find the best fit of that model to each other
  - o Typical methods: <span style="color:red">EM</span>, SOM

# Major Clustering Approaches (cont.)

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

❏ **Frequent Pattern-based approach**
  ○ Based on the analysis of frequent patterns
  ○ Typical methods: p-Cluster

❏ **Support Vector approach**
  ○ Based on the idea of mapping data points into higher dimensional feature space via a kernel function.
  ○ Typical methods: SVC, Kernel K-means

❏ **Graph Theoretic approach**
  ○ 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.
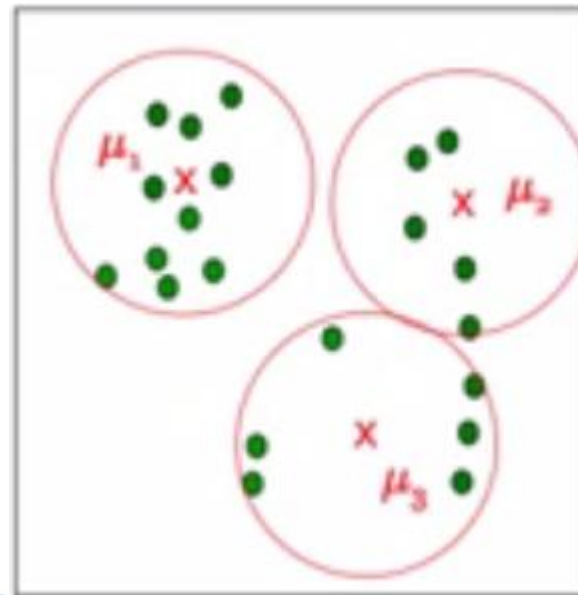
➢ Example: K-means

# K-means Clustering

❑ Assume $K$ clusters

❑ Iterate between two following steps:
  o Updating the assignment of data to clusters
  o Updating the cluster's summarization
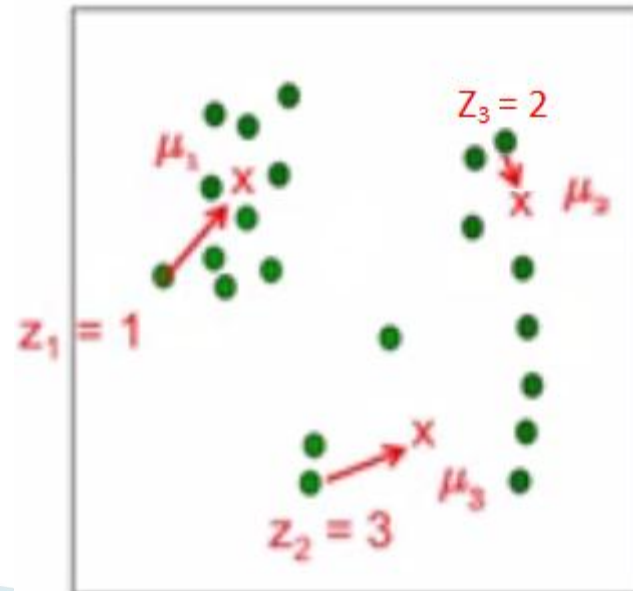  
  ———
  
  o sklearn.cluster.KMeans

# 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
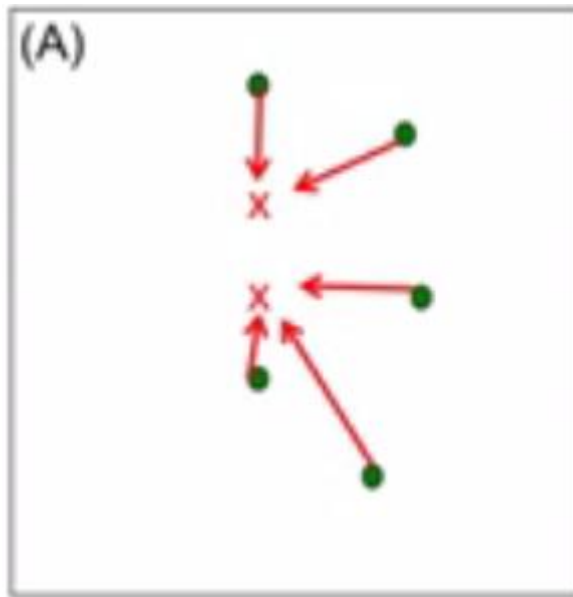- Assignment of $i_{th}$ example: $z_i \in 1..K$

# K-means Clustering (cont.)

❑ Iterate until convergence

   A. For each data, find the closest centroid:

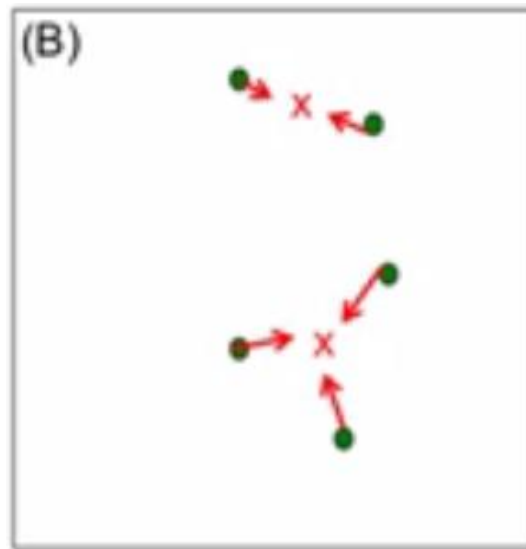$$z_i = \underset{c}{\text{argmin}} ||\boldsymbol{x_i} - \boldsymbol{\mu_c}||^2 , \forall i$$



(A)

# K-means Clustering (cont.)

❑ Iterate until convergence

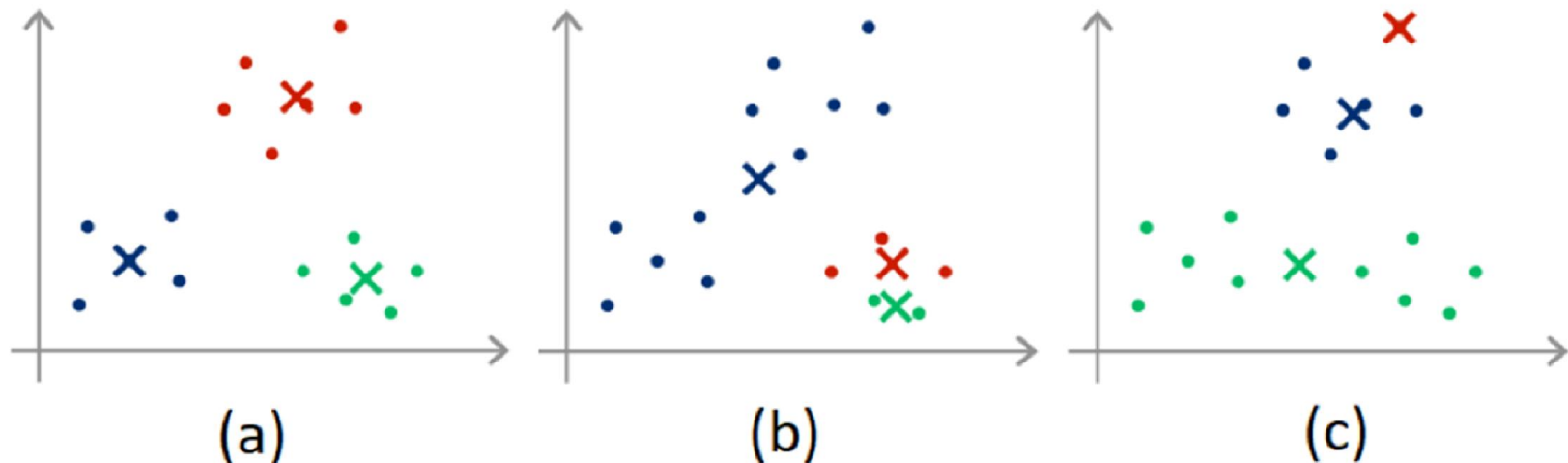B. Set each cluster to the mean of all assigned data:

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



(B)

Demo

# K-means Properties

☐ Poor initialization may lead to poor clustering



(a)　　　　　　(b)　　　　　　(c)

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

# K-means Properties (cont.)

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

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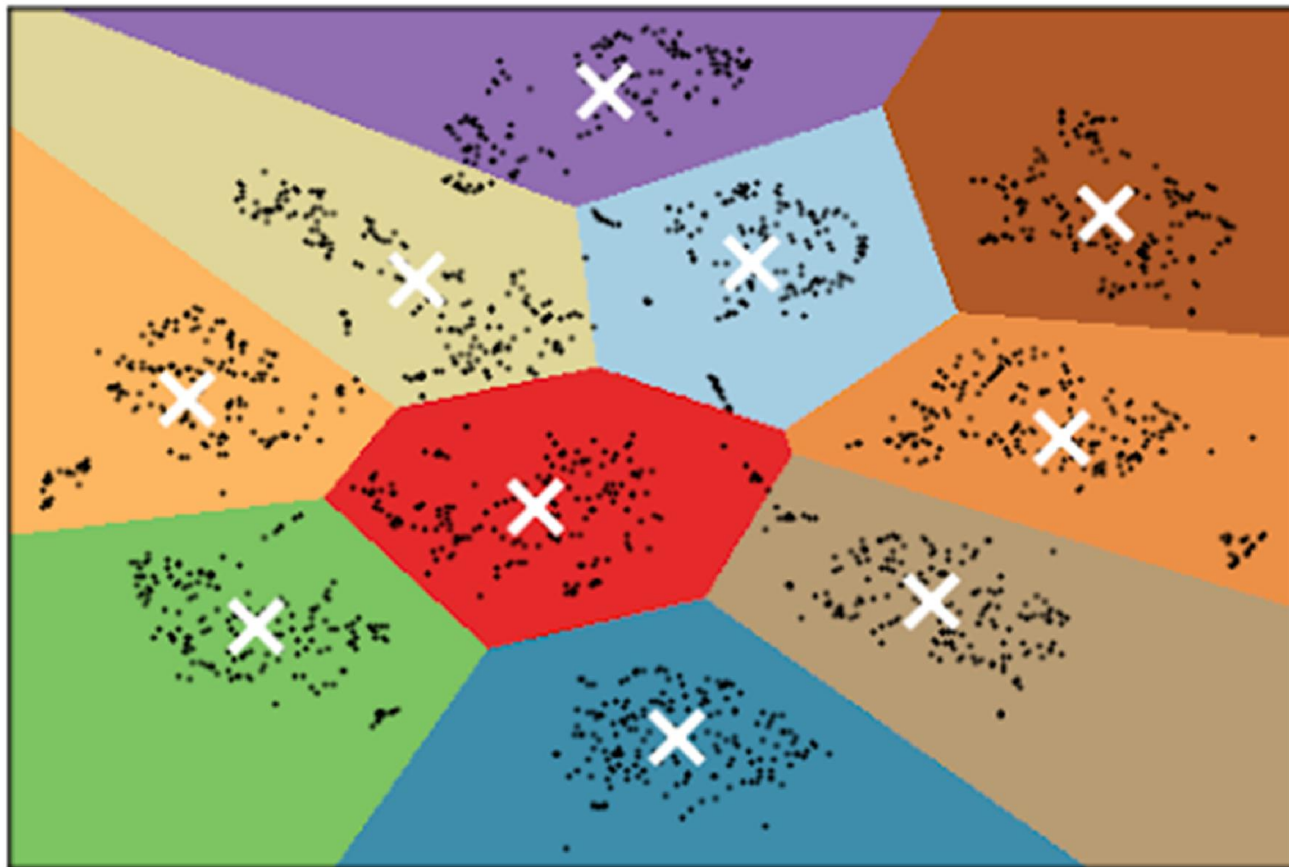    - Most commonly used on data when a mean or centroid cannot be defined, such as graphs.
  - …

# K-means Properties (cont.)

- Instance-based

- Time complexity: $O(tkm)$

- Non-parametric

- Linearly separable data

# K-means: Linear Separable
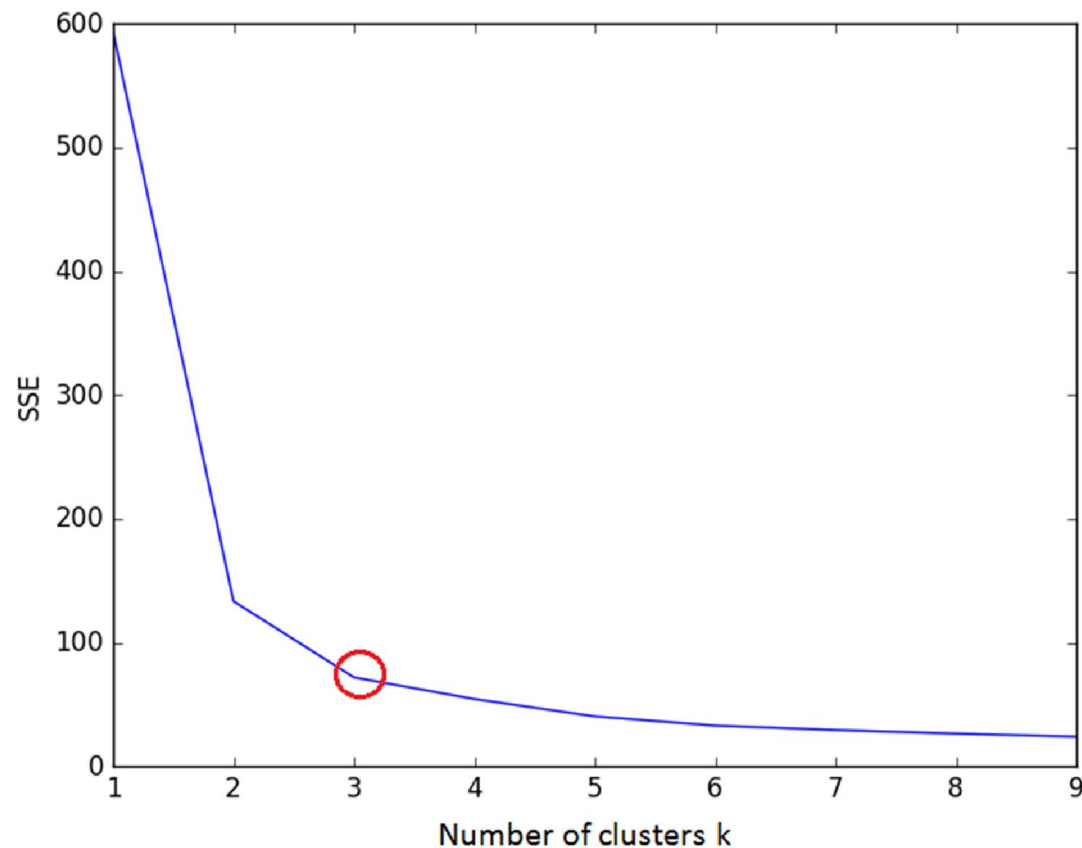
# Sum of Square Error

❑ Sum of Square Error (SSE)

$$SSE = \sum_k \sum_{x_i \in C_k} ||x_i - C_k||^2$$
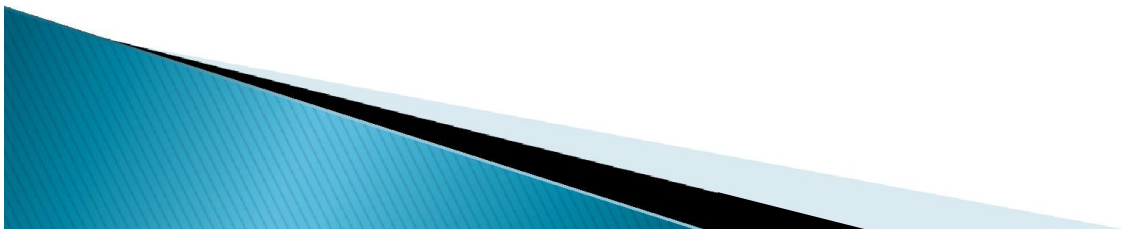
- Goal: minimizing within-cluster distance
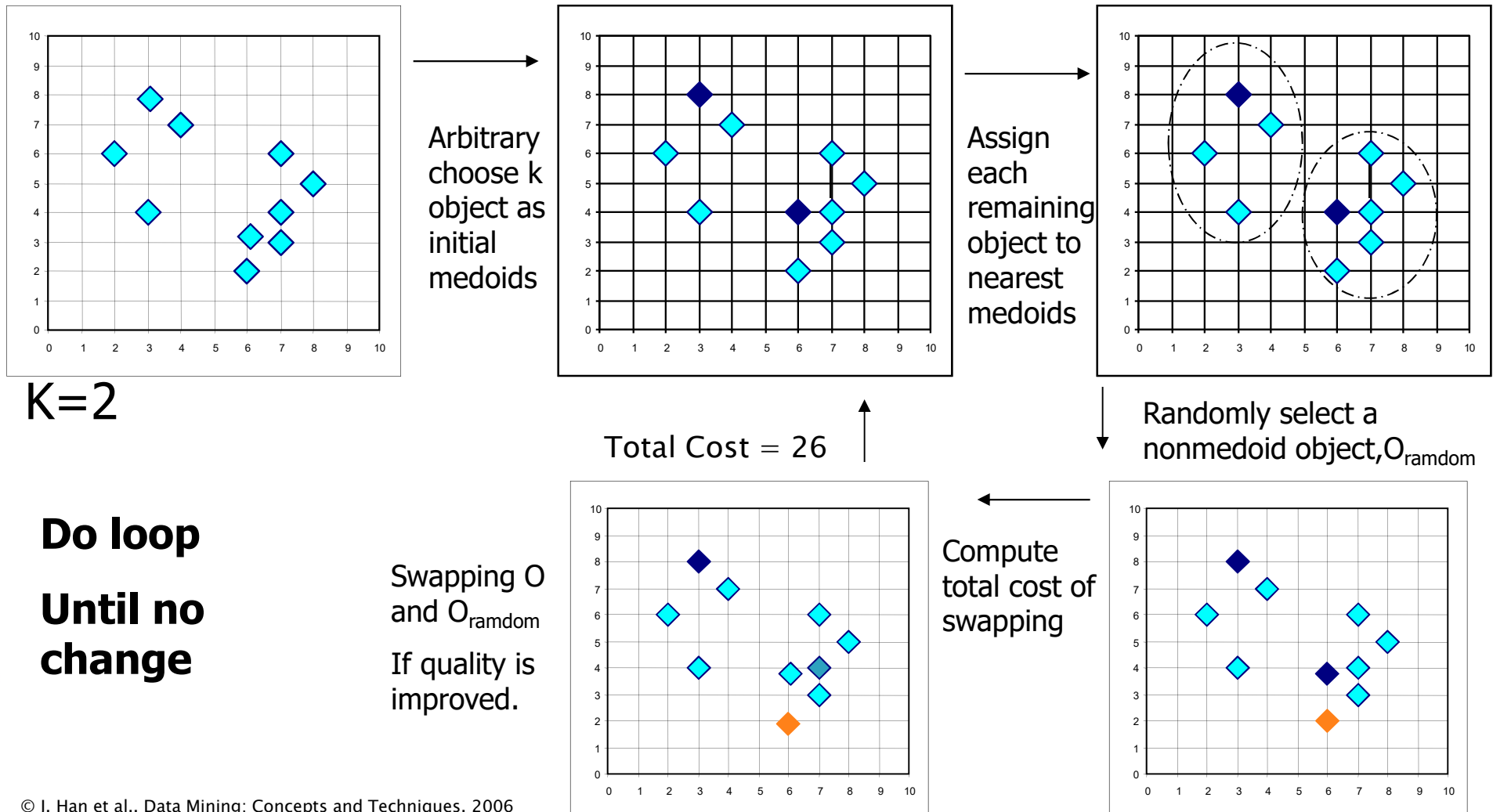
# Optimal number of Clusters

❑ Elbow method

# K-means Variations

- K-medoids/PAM (Partitioning Around Medoids)

- CLARA (Clustering Large Applications)

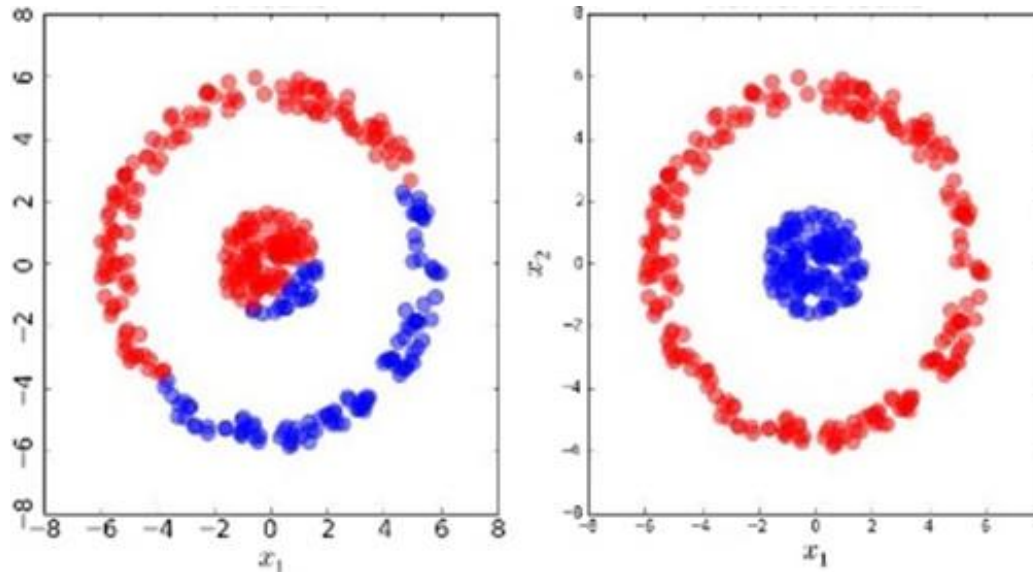- CLARANS (A Clustering Algorithm based on Randomized Search)

# PAM Algorithm

Total Cost = 20

Arbitrary choose k object as initial medoids

Assign each remaining object to nearest medoids

K=2

**Do loop**

**Until no change**

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

Total Cost = 26

Swapping O and $O_{ramdom}$

If quality is improved.

Compute total cost of swapping

# K-means Variations (cont.)

- Fuzzy C-means

- Kernel K-means

# Further Reading

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

□ Mean Shift Clustering

# References

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