




CS 412 Intro. to Data Mining

Chapter 10. Cluster Analysis: Basic Concepts and Methods

Jiawei Han, Computer Science, Univ. Illinois at Urbana-Champaign, 2017



Chapter 10. Cluster Analysis: Basic Concepts and Methods

- ☐ Cluster Analysis: An Introduction 
- ☐ Partitioning Methods
- ☐ Hierarchical Methods
- ☐ Density- and Grid-Based Methods
- ☐ Evaluation of Clustering (Coverage will be based on the available time)
- ☐ Summary

What Is Cluster Analysis?

- ❑ What is a cluster?
 - ❑ A cluster is a collection of data objects which are
 - ❑ Similar (or related) to one another within the same group (i.e., cluster)
 - ❑ Dissimilar (or unrelated) to the objects in other groups (i.e., clusters)
- ❑ Cluster analysis (or *clustering, data segmentation, ...*)
 - ❑ Given a set of data points, partition them into a set of groups (i.e., clusters) which are as similar as possible
- ❑ Cluster analysis is **unsupervised learning** (i.e., no predefined classes)
 - ❑ This contrasts with *classification* (i.e., *supervised learning*)
- ❑ Typical ways to use/apply cluster analysis
 - ❑ As a stand-alone tool to get insight into data distribution, or
 - ❑ As a preprocessing (or intermediate) step for other algorithms

What Is Good Clustering?

- ❑ A good clustering method will produce high quality clusters which should have
 - ❑ **High intra-class similarity:** **Cohesive** within clusters
 - ❑ **Low inter-class similarity:** **Distinctive** between clusters
- ❑ **Quality function**
 - ❑ 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
- ❑ There exist many similarity measures and/or functions for different applications
- ❑ Similarity measure is critical for cluster analysis

จำนวนจุดเริ่มต้นที่เลือก

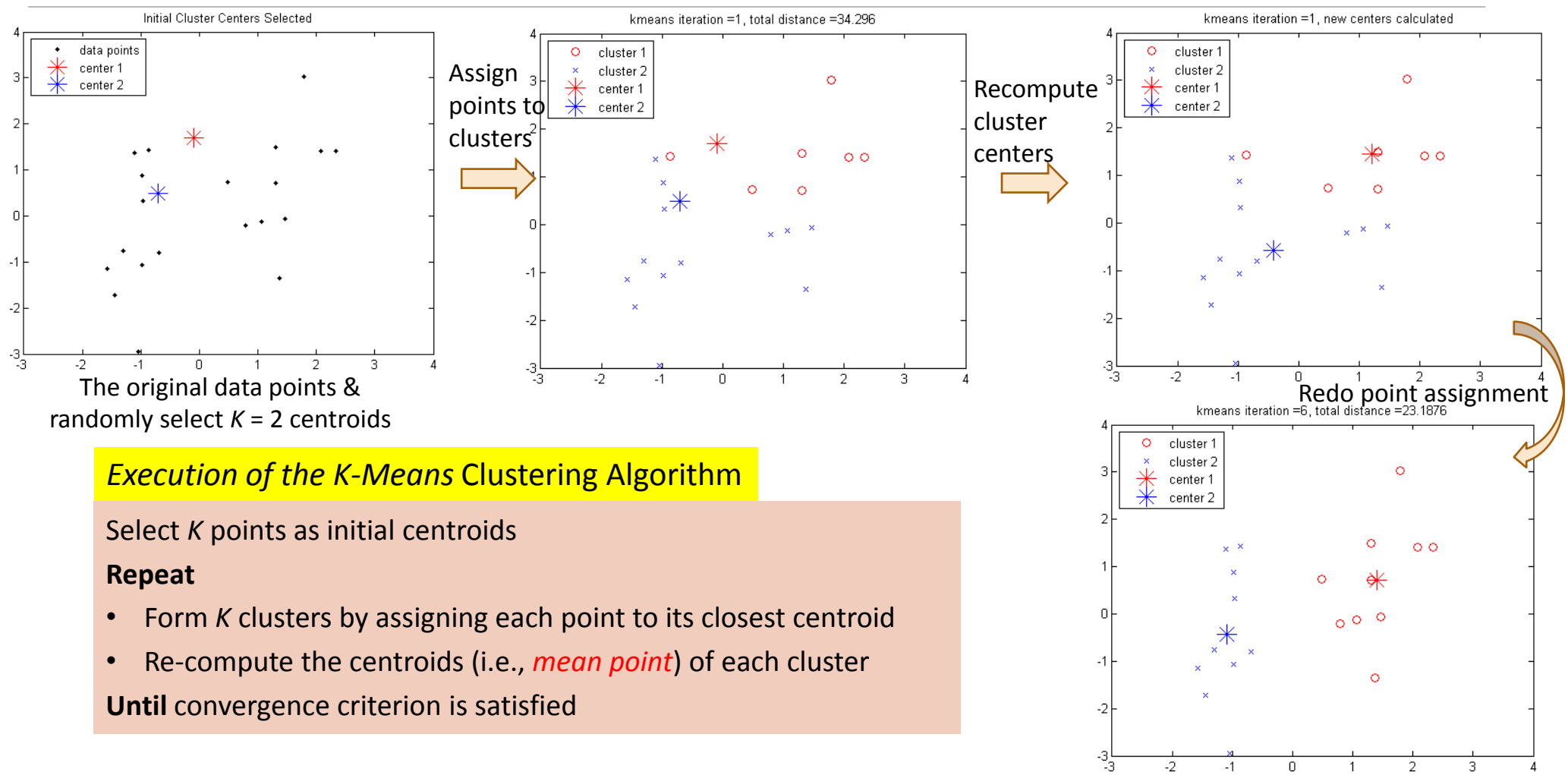
The *K-Means* Clustering Method

- *K-Means* (MacQueen'67, Lloyd'57/'82)
 - Each cluster is represented by the center of the cluster
- Given K , the number of clusters, the *K-Means* clustering algorithm is outlined as follows
 - Select K points as initial **centroids** จำนวนจุดเริ่มต้นที่เลือก
 - **Repeat** ทำซ้ำจนกว่าจะพอใจ
 - Form K clusters by assigning each point to its closest centroid
 - Re-compute the centroids (i.e., **mean point**) of each cluster
 - **Until** convergence criterion is satisfied
- Different kinds of measures can be used
 - Manhattan distance (L_1 norm), Euclidean distance (L_2 norm), Cosine similarity

เลือกจุด 2 จุด

จัดกลุ่มจุดที่ใกล้กว่าจุดใดจุดหนึ่ง

Example: *K-Means* Clustering



Variations of *K-Means*

- There are many variants of the *K-Means* method, varying in different aspects

- Choosing better initial centroid estimates

เลือกจุดเริ่มต้น

- K-means++*, *Intelligent K-Means*, *Genetic K-Means*

การเลือกจุดเริ่มต้นที่ดีขึ้น

To be discussed in this lecture

- Choosing different representative prototypes for the clusters

- K-Medoids*, *K-Medians*, *K-Modes*

K-means

To be discussed in this lecture

- Applying feature transformation techniques

- Weighted K-Means*, *Kernel K-Means*

การแปลงคุณลักษณะ

To be discussed in this lecture