



DATA MINING CLUSTERING

PRESENTED BY
DR. NEHA SHARMA, INDIA



2

OBJECTIVE

To understand-

- What is clustering?
- K-Means Cluster?



What is Cluster Analysis?

3

July 17, 2018

- Cluster: a collection of data objects
 - ▣ Similar to one another within the same cluster
 - ▣ Dissimilar to the objects in other clusters
- Cluster analysis
 - ▣ Finding similarities between data according to the characteristics found in the data and grouping similar data objects into clusters
- **Unsupervised learning**: no predefined classes
- Typical applications
 - ▣ As a **stand-alone tool** to get insight into data distribution
 - ▣ As a **preprocessing step** for other algorithms

Clustering: Rich Applications and Multidisciplinary Efforts

4

July 17, 2018

- Pattern Recognition
- Spatial Data Analysis
 - ▣ Create thematic maps in GIS by clustering feature spaces
 - ▣ Detect spatial clusters or for other spatial mining tasks
- Image Processing
- Economic Science (especially market research)
- WWW
 - ▣ Document classification
 - ▣ Cluster Weblog data to discover groups of similar access patterns

Examples of Clustering Applications

5

July 17, 2018

- Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs
- Land use: Identification of areas of similar land use in an earth observation database
- Insurance: Identifying groups of motor insurance policy holders with a high average claim cost
- 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

Quality: What Is Good Clustering?

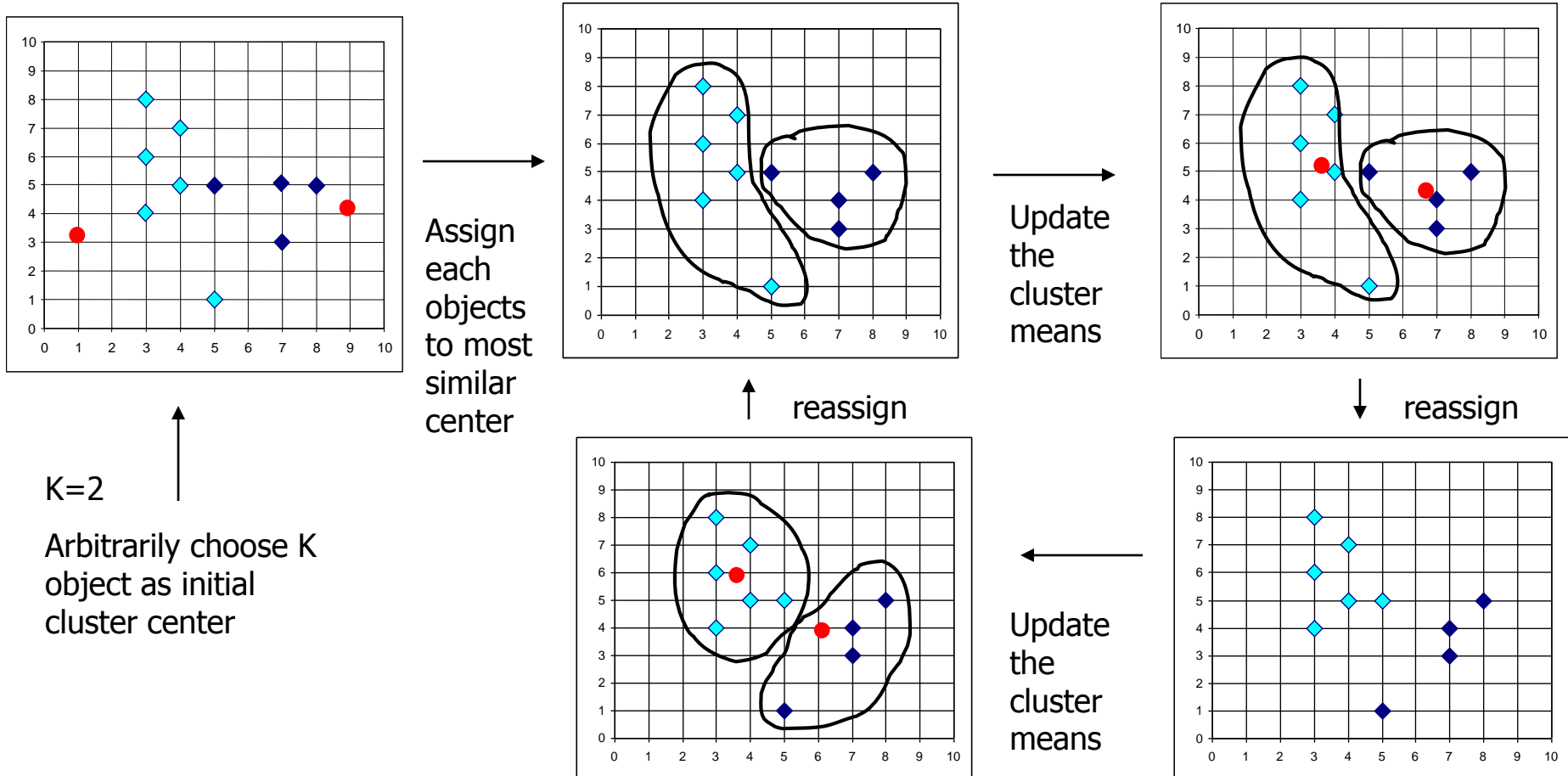
- A good clustering method will produce high quality clusters with
 - high intra-class similarity
 - low inter-class similarity
- The quality of a clustering result depends on both the similarity measure used by the method and its implementation
- The quality of a clustering method is also measured by its ability to discover some or all of the hidden patterns

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 partition (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 no more new assignment

The *K-Means* Clustering Method

■ Example



Comments on the *K-Means* Method

- Strength: *Relatively 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 optimum*. The *global optimum* may be found using techniques such as: *deterministic annealing* and *genetic algorithms*
- Weakness
 - Applicable only when *mean* is defined, then what about categorical data?
 - Need to specify k , the *number* of clusters, in advance
 - Unable to handle noisy data and *outliers*
 - Not suitable to discover clusters with *non-convex shapes*

Example - 1

10

July 17, 2018

Apply K-Means Clustering for following dataset for 2 clusters.
Tabulate the assignments.

Sample No	X	Y
1	185	72
2	170	56
3	168	60
4	179	68
5	182	72
6	188	77

Assignment - 1

11

July 17, 2018

Apply K-Means Clustering for following dataset for 2 clusters.
Tabulate the assignments.

	Driver_ID	Distance_Feature	Speeding_Feature
0	3423311935	71.24	28
1	3423313212	52.53	25
2	3423313724	64.54	27
3	3423311373	55.69	22
4	3423310999	54.58	25

12

Thank you

QUESTIONS????