

Clustering

Outline

- Cluster Analysis: Basic Concepts
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Evaluation of Clustering
- Summary

1. What is Clustering?

- **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
- **Clustering as Unsupervised learning:** no predefined classes (i.e., *learning by observations* vs. learning by examples: supervised)
- **Typical applications of clustering**
 - As a **stand-alone tool** to get insight into data distribution
 - As a **preprocessing step** for other algorithms

Examples in Security context

- Examples:
 - If someone (actor) is trying to *breach* your network -> likely they will try many times before actually getting through.
 - Or, if someone (actor) is sending *pharmaceutical spam*, they will need to send a lot of emails in order to get enough people to fall for the scam.
- => We can **segment** traffic into groups belonging to the same actor, and then block traffic from malicious actors.
- This process of segmentation is called **clustering**.

- Given a bunch of data points, which ones are similar to one another?
- For example, if you are trying to analyze a large dataset of *internet traffic to your site*, you might want to know which requests group together.
- => Some clusters might be *botnets*, some might be *mobile providers*, and some might be *legitimate users*.
- This process is called **clustering**.

Which clusters consist of malicious activity?

Clustering for Data Understanding and Applications

- **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

Clustering as a Preprocessing Tool (Utility)

- **Summarization:**
 - Preprocessing for regression, PCA, classification, and association analysis
- **Compression:**
 - Image processing: vector quantization
- **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

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

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)

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

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

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

2. 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)

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

(Loss function)

- **In more detailed**: 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*: Each cluster is represented by the center of the cluster
 - *k-medoids* or PAM (Partition around medoids): Each cluster is represented by one of the objects in the cluster

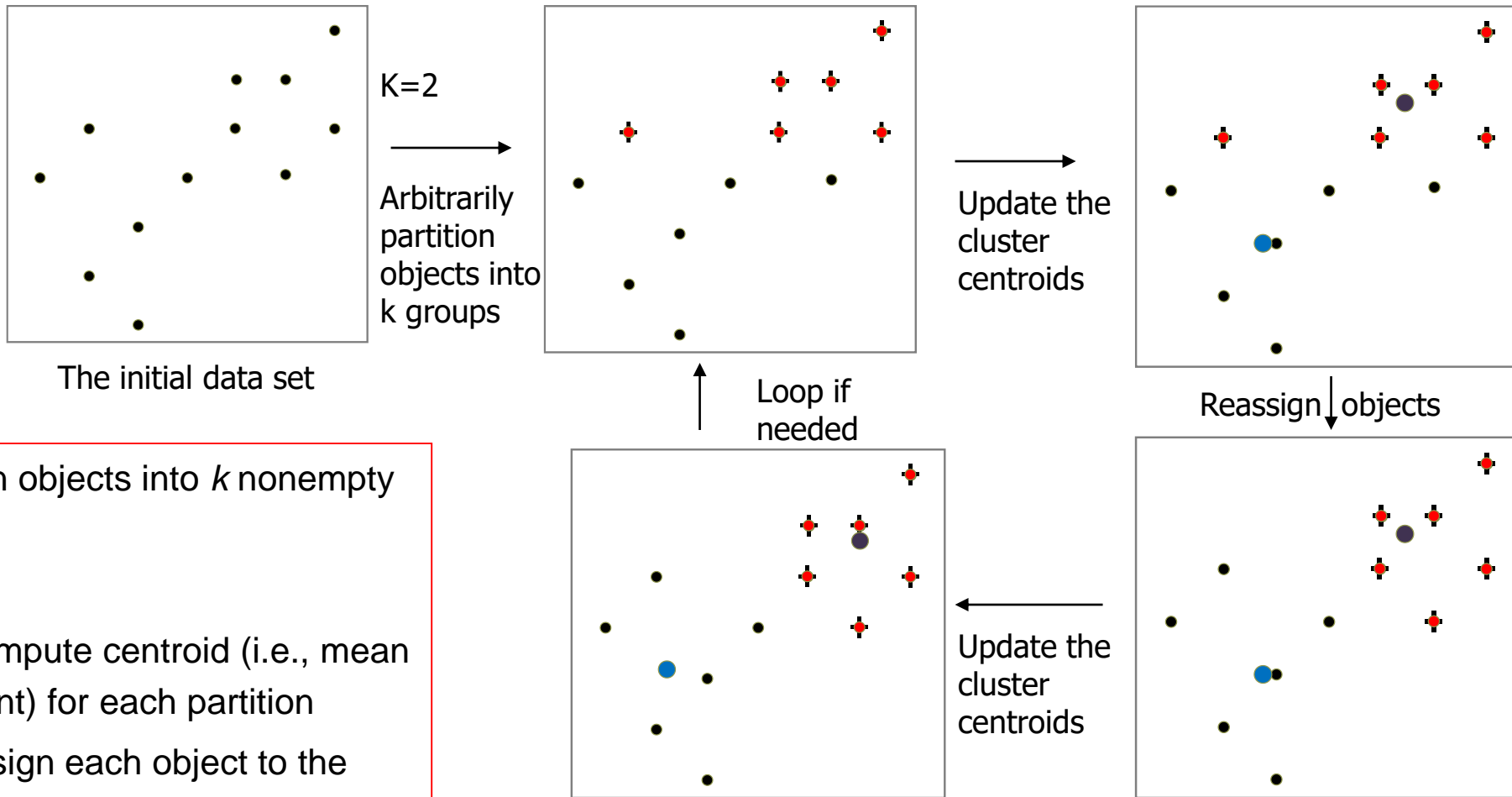
Mathematically, a loss function is a function that maps a set of pairs of (predicted label, truth label) to a real number.

The goal of a machine learning algorithm is to find the model parameters that produce predicted labels for the training set that minimize the loss function

The *K-Means* Clustering Method

- Given k , the *k-means* algorithm is implemented in four steps:
 1. Partition objects into k nonempty subsets
 2. 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)
 3. Assign each object to the cluster with the nearest seed point
 4. Go back to Step 2, stop when the assignment does not change (the difference between loss values on successive iterations is below a predetermined threshold)

An Example of *K-Means* Clustering

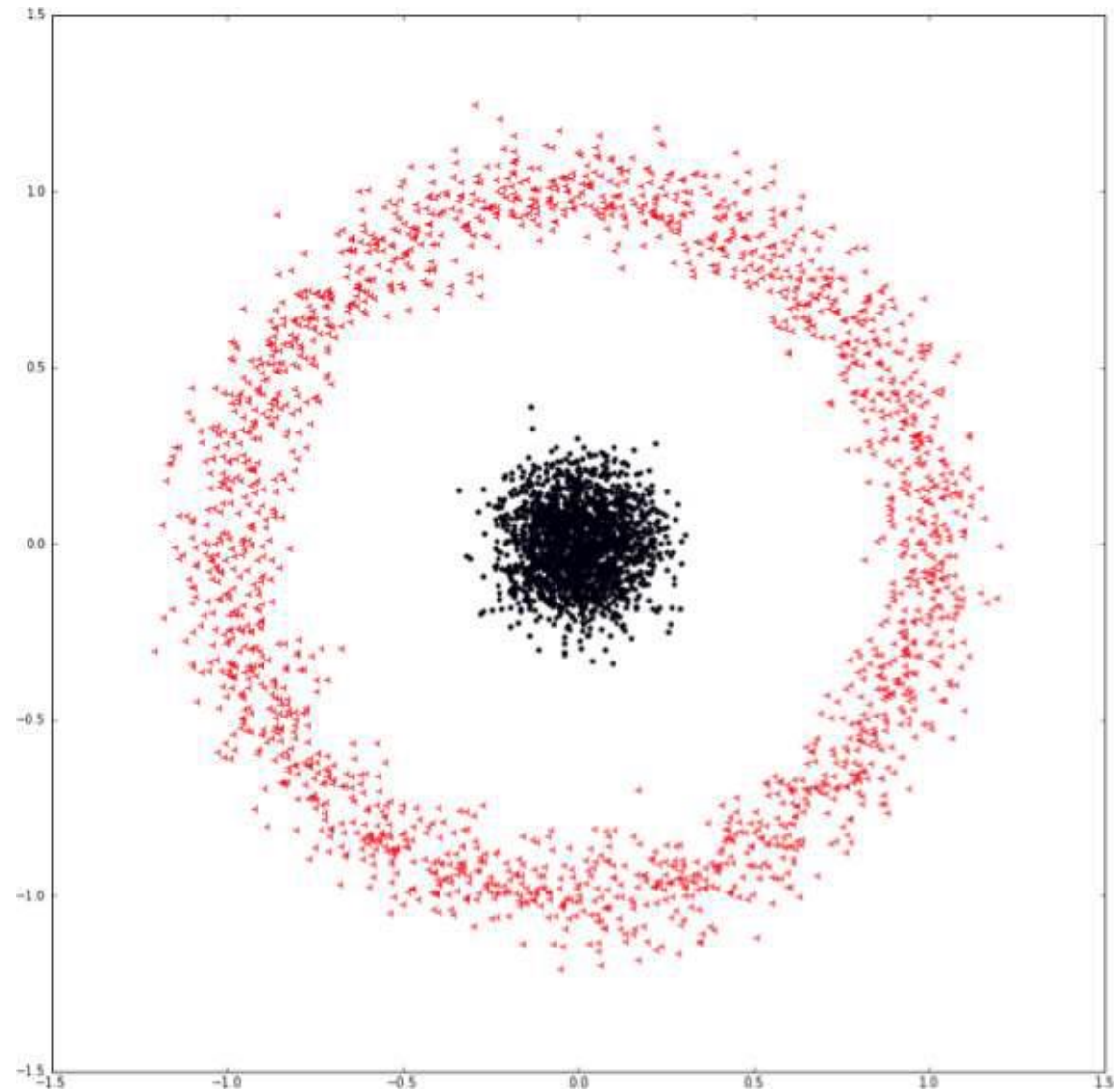


- Partition objects into k nonempty subsets
- Repeat
 - Compute centroid (i.e., mean point) for each partition
 - Assign each object to the cluster of its nearest centroid
- Until no change

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
 - Need to specify k , the *number* of clusters, in advance.
 - Sensitive to noisy data and *outliers*
 - Not suitable to discover clusters with *non-convex shapes*

What if?

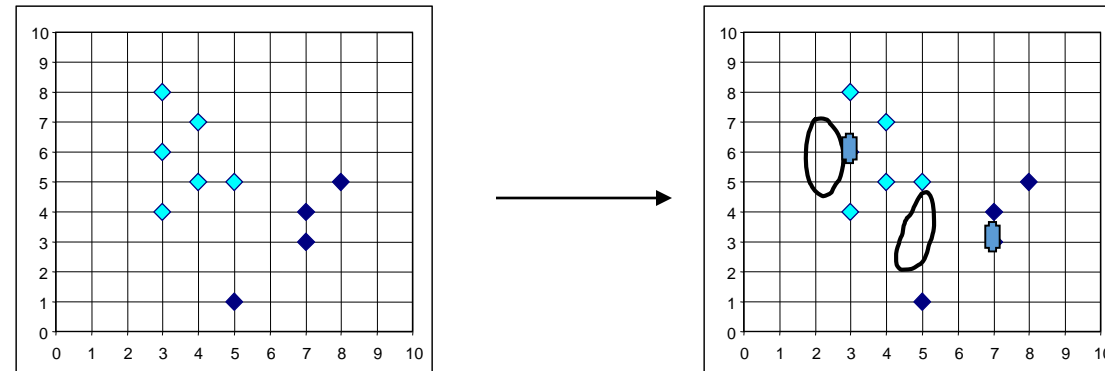


Variations of the *K-Means* Method

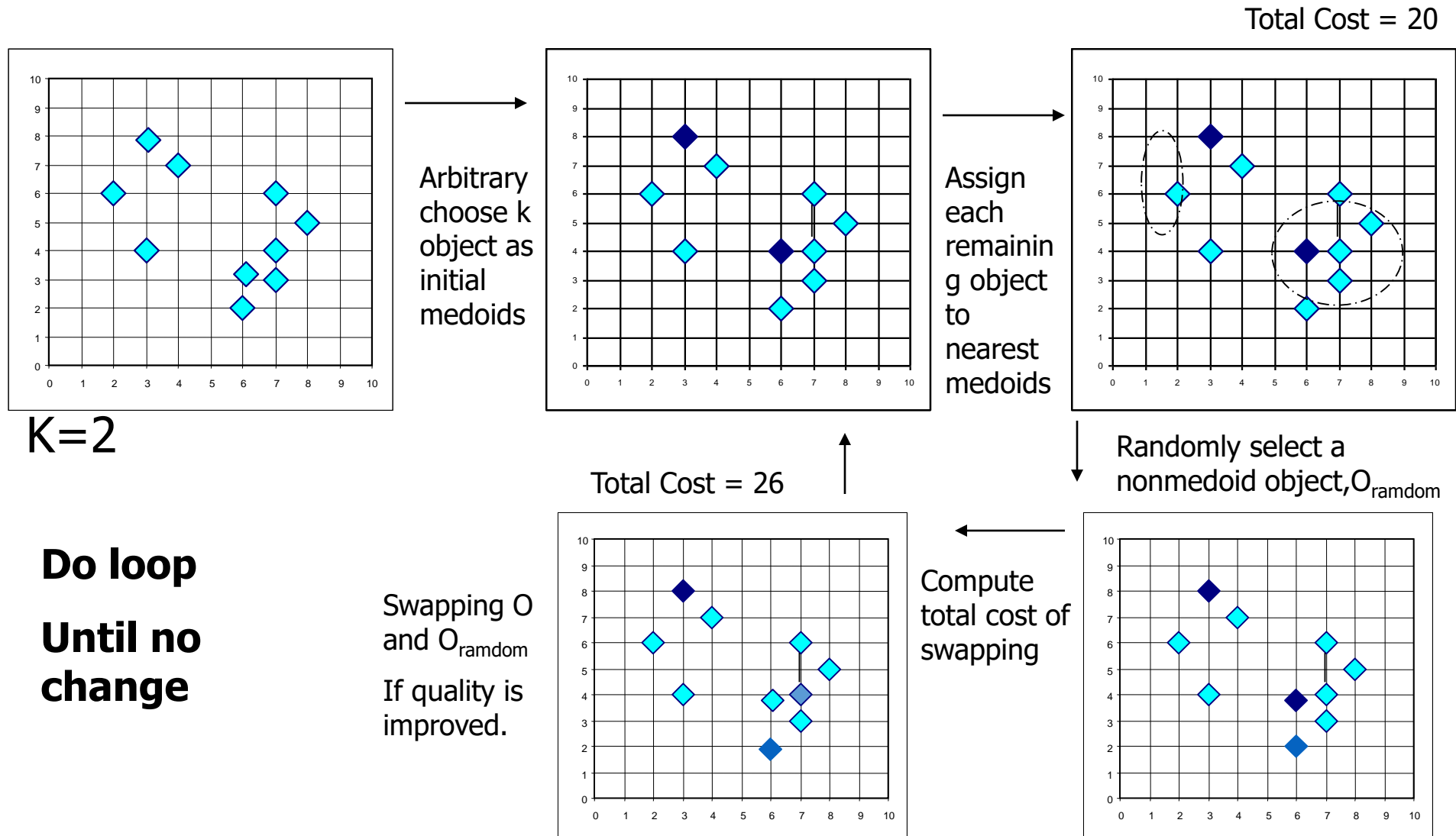
- Most of the variants of the *k-means* which differ in
 - Selection of the initial *k* means
 - Dissimilarity calculations
 - Strategies to calculate cluster means
- Handling categorical data: *k-modes*
 - Replacing means of clusters with modes (*the most common*)
 - Using new dissimilarity measures 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

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



PAM: A Typical K-Medoids Algorithm

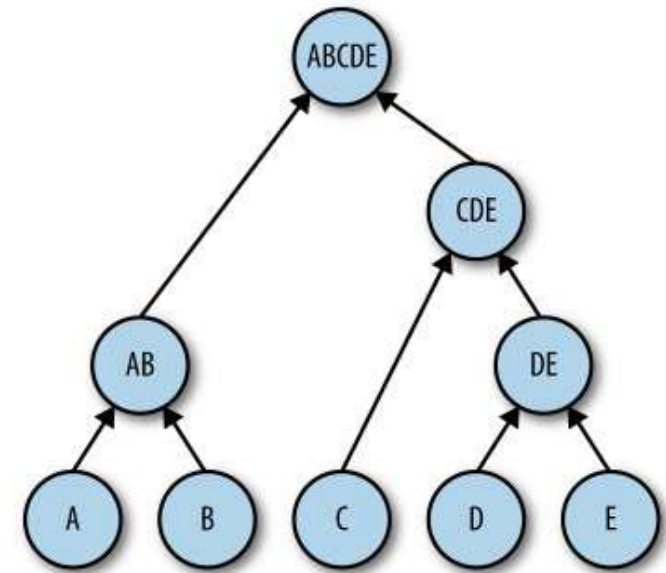


The K-Medoid Clustering Method

- *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
 - *CLARA* (Kaufmann & Rousseeuw, 1990): PAM on samples
 - *CLARANS* (Ng & Han, 1994): Randomized re-sampling

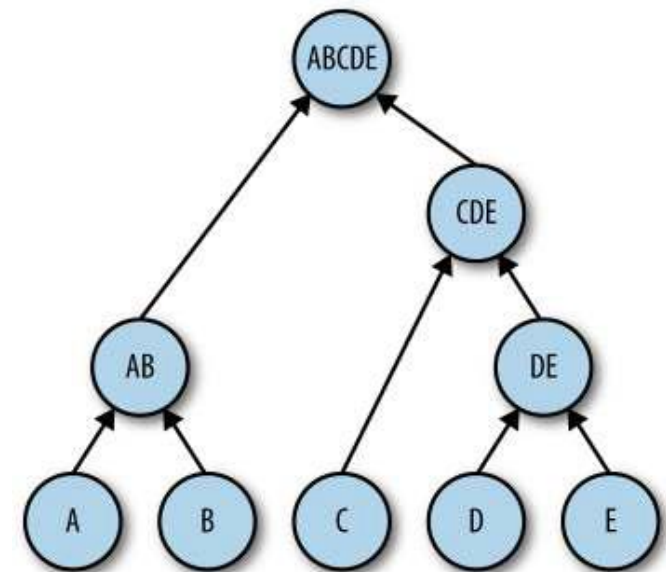
3. Hierarchical clustering

- Unlike the *k-means* algorithm, hierarchical clustering methods are not parametrized by an operator-selected value k (the number of clusters you want to create).
- Choosing an appropriate k is a nontrivial task, and can significantly affect clustering results.



3. Hierarchical clustering

- Agglomerative (bottom-up) hierarchical clustering builds clusters as follows
 - 1. Assign each data point to its own cluster (see Figure, bottom layer).
 - 2. Merge the two clusters that are the most similar, where “most similar” is determined by a distance metric such as the Euclidean distance.
 - 3. Repeat step 2 until there is only one cluster remaining (see Figure, top layer).
 - 4. Navigate the layers of this tree (dendrogram) and select the layer that gives you the most appropriate clustering result.



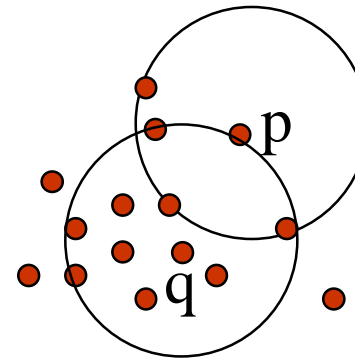
4. Density-Based Clustering Methods

- Clustering based on density (local cluster criterion), such as density-connected points
- Major features:
 - Discover clusters of arbitrary shape
 - Handle noise
 - One scan
 - Need density parameters as termination condition
- Several interesting studies:
 - DBSCAN: Ester, et al. (KDD'96)
 - OPTICS: Ankerst, et al (SIGMOD'99).
 - DENCLUE: Hinneburg & D. Keim (KDD'98)
 - CLIQUE: Agrawal, et al. (SIGMOD'98) (more grid-based)

Density-Based Clustering: Basic Concepts

- Two parameters:
 - *Eps*: Maximum radius of the neighbourhood
 - *MinPts*: Minimum number of points in an *Eps*-neighbourhood of that point
- $N_{Eps}(p)$: $\{q \text{ belongs to } D \mid \text{dist}(p,q) \leq Eps\}$
- **Directly density-reachable**: A point p is directly density-reachable from a point q w.r.t. *Eps*, *MinPts* if
 - p belongs to $N_{Eps}(q)$
 - q satisfies core point condition:

$$|N_{Eps}(q)| \geq MinPts$$



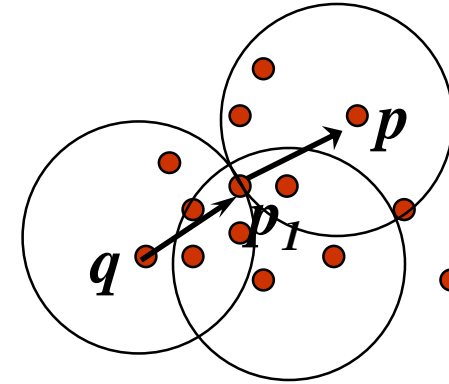
MinPts = 5

Eps = 1 cm

Density-Reachable and Density-Connected

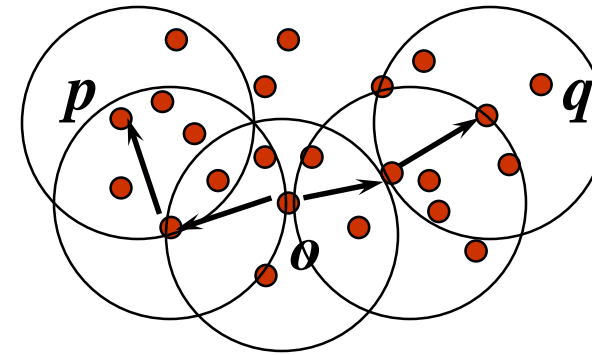
- Density-reachable:

- A point p is **density-reachable** from a point q w.r.t. Eps , $MinPts$ if there is a chain of points p_1, \dots, p_n , $p_1 = q$, $p_n = p$ such that p_{i+1} is directly density-reachable from p_i



- Density-connected

- A point p is **density-connected** to a point q w.r.t. Eps , $MinPts$ if there is a point o such that both, p and q are density-reachable from o w.r.t. Eps and $MinPts$



DBSCAN: Density-Based Spatial Clustering of Applications with Noise

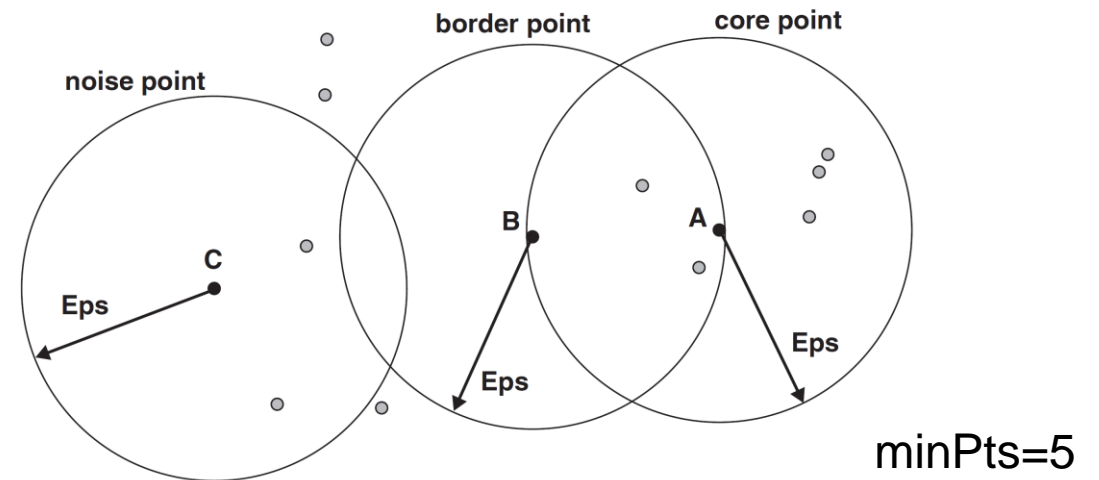
- Relies on a *density-based* notion of cluster: A *cluster* is defined as a maximal set of density-connected points
- Discovers clusters of arbitrary shape in spatial databases with noise

Each data point is classified as a core point, a border point, or a noise point:

—Core points are points that have at least minPts number of points within their ϵ -radius.

—Border points are themselves not core points, but are covered within the ϵ -radius of some core point.

—Noise points (Outliers) are neither core points nor border points.



https://www2.cs.uh.edu/~ceick/UDM/dm_clustering2.pptx

DBSCAN: The Algorithm

1. Arbitrary select a point p (mark it *visited*)
2. Retrieve all points density-reachable from p w.r.t. Eps and $MinPts$
3. If p is a core point, a cluster is formed
4. Continue the process until all of the points have been processed

5. Determine the Number of Clusters

- Empirical method
 - # of clusters $\approx \sqrt{n}/2$ for a dataset of n points
- Cross validation method
 - Divide a given data set into m parts
 - Use $m - 1$ parts to obtain a clustering model
 - Use the remaining part to test the quality of the clustering
 - E.g., For each point in the test set, find the closest centroid, and use the sum of squared distance between all points in the test set and the closest centroids to measure how well the model fits the test set
 - For any $k > 0$, repeat it m times, compare the overall quality measure w.r.t. different k 's, and find # of clusters that fits the data the best

Measuring Clustering Quality

- Two methods: extrinsic vs. intrinsic
- Extrinsic: supervised, i.e., the ground truth is available
 - Compare a clustering against the ground truth using certain clustering quality measure
 - Ex. BCubed precision and recall metrics
- Intrinsic: unsupervised, i.e., the ground truth is unavailable
 - Evaluate the goodness of a clustering by considering how well the clusters are separated, and how compact the clusters are
 - Ex. Silhouette coefficient

Silhouette coefficient

This score is calculated separately for each sample in the dataset. Using some distance metric (e.g., Euclidean distance), we find the following two mean distances for some sample x :

- a : the mean distance between sample x and all other samples in the same cluster
- b : the mean distance between sample x and all other samples in the next nearest cluster

The Silhouette coefficient s is then defined as follows:⁴⁰

$$s = \frac{b - a}{\max(a, b)}$$

Measuring Clustering Quality: Extrinsic Methods

- Clustering quality measure: $Q(C, C_g)$, for a clustering C given the ground truth C_g .
- Q is good if it satisfies the following **4** essential criteria
 - Cluster homogeneity: the purer, the better
 - Cluster completeness: should assign objects belong to the same category in the ground truth to the same cluster
 - Rag bag: putting a heterogeneous object into a pure cluster should be penalized more than putting it into a *rag bag* (i.e., “miscellaneous” or “other” category)
 - Small cluster preservation: splitting a small category into pieces is more harmful than splitting a large category into pieces

Summary

- **Cluster analysis** groups objects based on their **similarity** and has wide applications
- Measure of similarity can be computed for **various types of data**
- Clustering algorithms can be **categorized** into partitioning methods, hierarchical methods, density-based methods, grid-based methods, and model-based methods
- Quality of clustering results can be evaluated in various ways