

Data Mining -- Clustering

Instructor: Jen-Wei Huang

Office: 92528 in the EE building
jwhuang@mail.ncku

Clustering

- ▶ 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 or as a **preprocessing step** for other algorithms
- ▶ Pattern Recognition
- ▶ Image Processing
- ▶ Economic Science (especially market research)
- ▶ WWW
 - Web pages (resources) clustering
 - Cluster Weblog data to discover groups of similar access patterns

Examples

- ▶ 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

Quality of Clustering

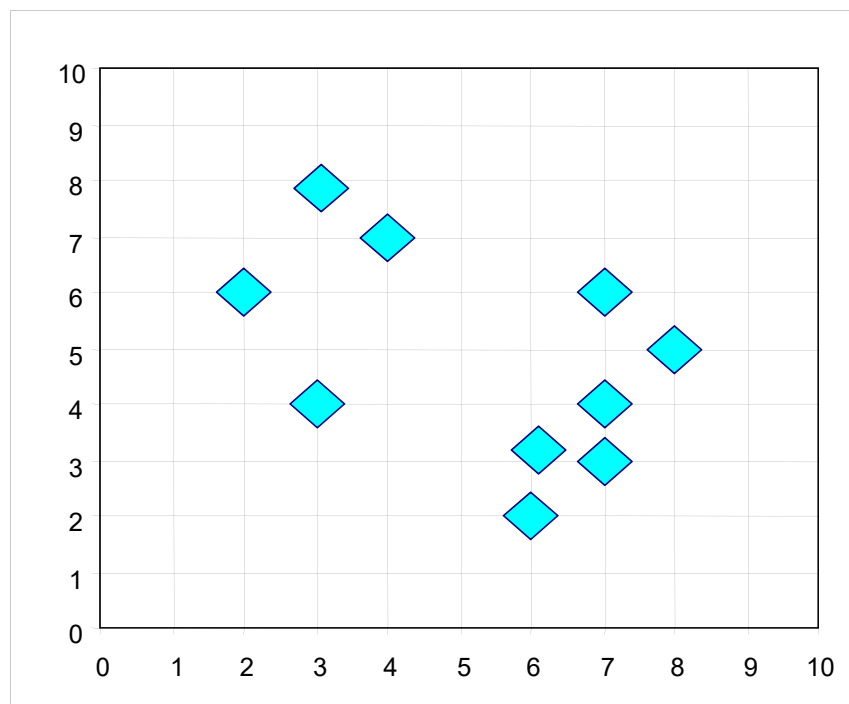
- ▶ Good clusters :
 - high intra-class similarity
 - low inter-class similarity
- ▶ The quality of a clustering result depends on both the similarity measure used by the method and its implementation
- ▶ The quality of a clustering method is also measured by its ability to discover some or all of the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

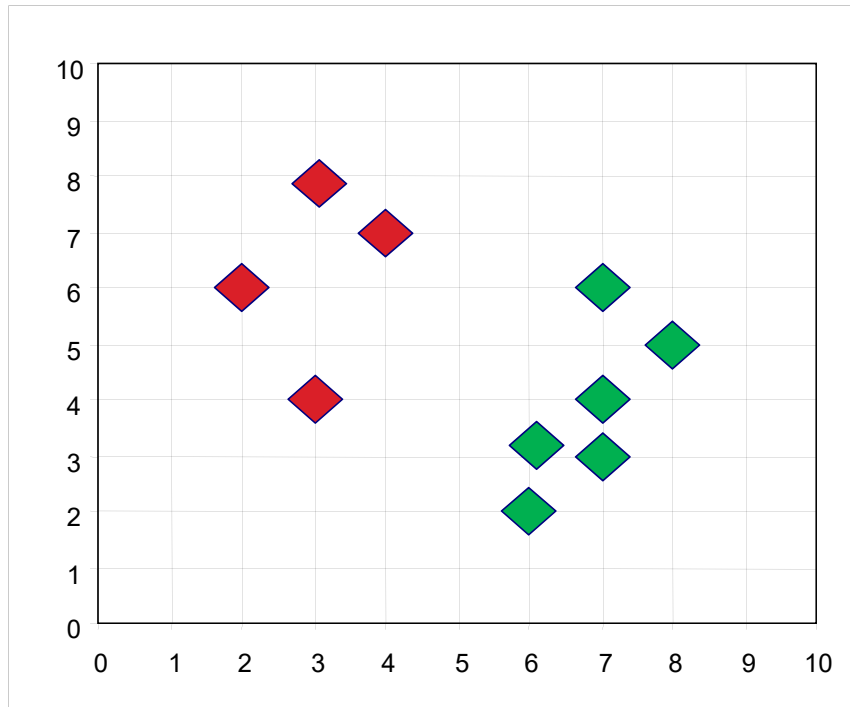
Measure the Quality

- ▶ **Dissimilarity/Similarity metric**: Similarity is expressed in terms of a distance function, typically metric: $d(i, j)$
- ▶ There is a separate “quality” function that measures the “goodness” of a cluster.
- ▶ The definitions of **distance functions** are usually very different for interval-scaled, boolean, categorical, ordinal ratio, and vector variables.
- ▶ Weights should be associated with different variables based on applications and data semantics.
- ▶ It is hard to define “similar enough” or “good enough”
 - the answer is typically highly subjective.

Requirements of Clustering

- ▶ Able to deal with noise and outliers
- ▶ Insensitive to order of input records
- ▶ High dimensionality
- ▶ Incorporation of user-specified constraints
- ▶ Interpretability and usability
- ▶ Scalability
- ▶ Ability to deal with different types of attributes
- ▶ Ability to handle dynamic data
- ▶ Discovery of clusters with arbitrary shape
- ▶ Determination input parameters





Type of Data

- ▶ Interval-scaled variables
- ▶ Binary variables
- ▶ Nominal, ordinal, and ratio variables
- ▶ Variables of mixed types

Interval-Valued Variables

- ▶ Standardize data

- Calculate the mean absolute deviation:

$$s_f = \frac{1}{n}(|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|)$$

where $m_f = \frac{1}{n}(x_{1f} + x_{2f} + \dots + x_{nf})$.

- Calculate the standardized measurement (z-score)

$$z_{if} = \frac{x_{if} - m_f}{s_f}$$

- ▶ Using mean absolute deviation is more robust than using standard deviation

Distance

- ▶ Distances are normally used to measure the similarity or dissimilarity between two data objects

- ▶ Some popular ones include: *Minkowski distance*:

$$d(i, j) = \sqrt[q]{(|x_{i1} - x_{j1}|^q + |x_{i2} - x_{j2}|^q + \dots + |x_{ip} - x_{jp}|^q)}$$

where $i = (x_{i1}, x_{i2}, \dots, x_{ip})$ and $j = (x_{j1}, x_{j2}, \dots, x_{jp})$ are two p -dimensional data objects, and q is a positive integer

- ▶ If $q = 1$, d is Manhattan distance

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

Distance

- ▶ If $q = 2$, d is Euclidean distance:

$$d(i, j) = \sqrt{(|x_{i_1} - x_{j_1}|^2 + |x_{i_2} - x_{j_2}|^2 + \dots + |x_{i_p} - x_{j_p}|^2)}$$

- Properties

- $d(i, j) \geq 0$
- $d(i, i) = 0$
- $d(i, j) = d(j, i)$
- $d(i, j) \leq d(i, k) + d(k, j)$

- ▶ Also, one can use weighted distance, parametric Pearson product moment correlation, or other dissimilarity measures

Binary Variables

- ▶ A contingency table for binary data

		Object j		
		1	0	sum
Object i	1	a	b	$a+b$
	0	c	d	$c+d$
sum		$a+c$	$b+d$	p

- ▶ Distance measure for symmetric binary variables:
- ▶ Distance measure for asymmetric binary variables:
- ▶ Jaccard coefficient (*similarity* measure for *asymmetric* binary variables):

$$d(i, j) = \frac{b + c}{a + b + c + d}$$

$$d(i, j) = \frac{b + c}{a + b + c}$$

$$sim_{Jaccard}(i, j) = \frac{a}{a + b + c}$$

Dissimilarity

▶ Example

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

- gender is a symmetric attribute
- the remaining attributes are asymmetric binary
- let the values Y and P be set to 1, and the value N be set to 0

$$d(\text{jack}, \text{mary}) = \frac{0 + 1}{2 + 0 + 1} = 0.33$$

$$d(\text{jack}, \text{jim}) = \frac{1 + 1}{1 + 1 + 1} = 0.67$$

$$d(\text{jim}, \text{mary}) = \frac{1 + 2}{1 + 1 + 2} = 0.75$$

Nominal Variables

- ▶ A generalization of the binary variable in that it can take more than 2 states, e.g., red, yellow, blue, green
- ▶ Method 1: Simple matching
 - m : # of matches, p : total # of variables

$$d(i, j) = \frac{p - m}{p}$$

- ▶ Method 2: use a large number of binary variables
 - creating a new binary variable for each of the M nominal states

Ordinal Variables

- ▶ An ordinal variable can be discrete or continuous
- ▶ Order is important, e.g., rank
- ▶ Can be treated like interval-scaled
 - replace x_{if} by their rank $r_{if} \in \{1, \dots, M_f\}$
 - map the range of each variable onto $[0, 1]$ by replacing i -th object in the f -th variable by

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

- compute the dissimilarity using methods for interval-scaled variables

Ratio-Scaled Variables

- ▶ Ratio-scaled variable: a positive measurement on a nonlinear scale, approximately at exponential scale, such as Ae^{Bt} or Ae^{-Bt}
- ▶ Methods:
 - treat them like interval-scaled variables—*not a good choice!* (why?—the scale can be distorted)
 - apply logarithmic transformation

$$y_{if} = \log(x_{if})$$

- treat them as continuous ordinal data treat their rank as interval-scaled

Variables of Mixed Types

- ▶ A database may contain all the six types of variables
 - symmetric binary, asymmetric binary, nominal, ordinal, interval and ratio
- ▶ One may use a weighted formula to combine their effects

$$d(i, j) = \frac{\sum_{f=1}^p \delta_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^p \delta_{ij}^{(f)}}$$

- f is binary or nominal: $d_{ij}^{(f)} = 0$ if $x_{if} = x_{jf}$,
or $d_{ij}^{(f)} = 1$ otherwise
- f is interval-based: use the normalized distance
- f is ordinal or ratio-scaled
 - compute ranks r_{if} and
 - and treat z_{if} as interval-scaled $z_{if} = \frac{r_{if} - 1}{M_f - 1}$

Vector Objects

- ▶ Vector objects: keywords in documents, gene features in micro-arrays, etc.
- ▶ Broad applications: information retrieval, biologic taxonomy, etc.
- ▶ Cosine measure

$$s(\vec{X}, \vec{Y}) = \frac{\vec{X}^t \cdot \vec{Y}}{|\vec{X}| |\vec{Y}|}$$

- ▶ A variant: Tanimoto coefficient

$$s(\vec{X}, \vec{Y}) = \frac{\vec{X}^t \cdot \vec{Y}}{\vec{X}^t \cdot \vec{X} + \vec{Y}^t \cdot \vec{Y} - \vec{X}^t \cdot \vec{Y}}$$

Clustering Methods

- ▶ 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, ROCK, CAMELEON
- ▶ Density-based approach:
 - Based on connectivity and density functions
 - Typical methods: DBSACN, OPTICS, DenClue

Clustering Methods

- ▶ Grid-based approach:
 - based on a multiple-level granularity structure
 - Typical methods: STING, WaveCluster, CLIQUE
- ▶ 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: pCluster

Clustering Methods

- ▶ User-guided or constraint-based:
 - Clustering by considering user-specified or application-specific constraints
 - Typical methods: COD (obstacles), constrained clustering

Distance between Clusters

- ▶ Single link: smallest distance between an element in one cluster and an element in the other, i.e., $\text{dis}(K_i, K_j) = \min(t_{ip}, t_{jq})$
- ▶ Complete link: largest distance between an element in one cluster and an element in the other, i.e., $\text{dis}(K_i, K_j) = \max(t_{ip}, t_{jq})$
- ▶ Average: avg distance between an element in one cluster and an element in the other, i.e., $\text{dis}(K_i, K_j) = \text{avg}(t_{ip}, t_{jq})$
- ▶ Centroid: distance between the centroids of two clusters, i.e., $\text{dis}(K_i, K_j) = \text{dis}(C_i, C_j)$
- ▶ Medoid: distance between the medoids of two clusters, i.e., $\text{dis}(K_i, K_j) = \text{dis}(M_i, M_j)$
 - Medoid: one chosen, centrally located object in the cluster

Centroid, Radius and Diameter

- ▶ Centroid: the “middle” of a cluster

$$C_m = \frac{\sum_{i=1}^N (t_{ip})}{N}$$

- ▶ Radius: square root of average distance from any point of the cluster to its centroid

$$R_m = \sqrt{\frac{\sum_{i=1}^N (t_{ip} - c_m)^2}{N}}$$

- ▶ Diameter: square root of average mean squared distance between all pairs of points in the cluster

$$D_m = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N (t_{ip} - t_{jq})^2}{N(N-1)}}$$

Partitioning Algorithms: Basic Concept

- ▶ Partitioning method: Construct a partition of a database D of n objects into a set of k clusters, s.t., min sum of squared distance

$$\sum_{m=1}^k \sum_{t_{mi} \in K_m} (C_m - t_{mi})^2$$

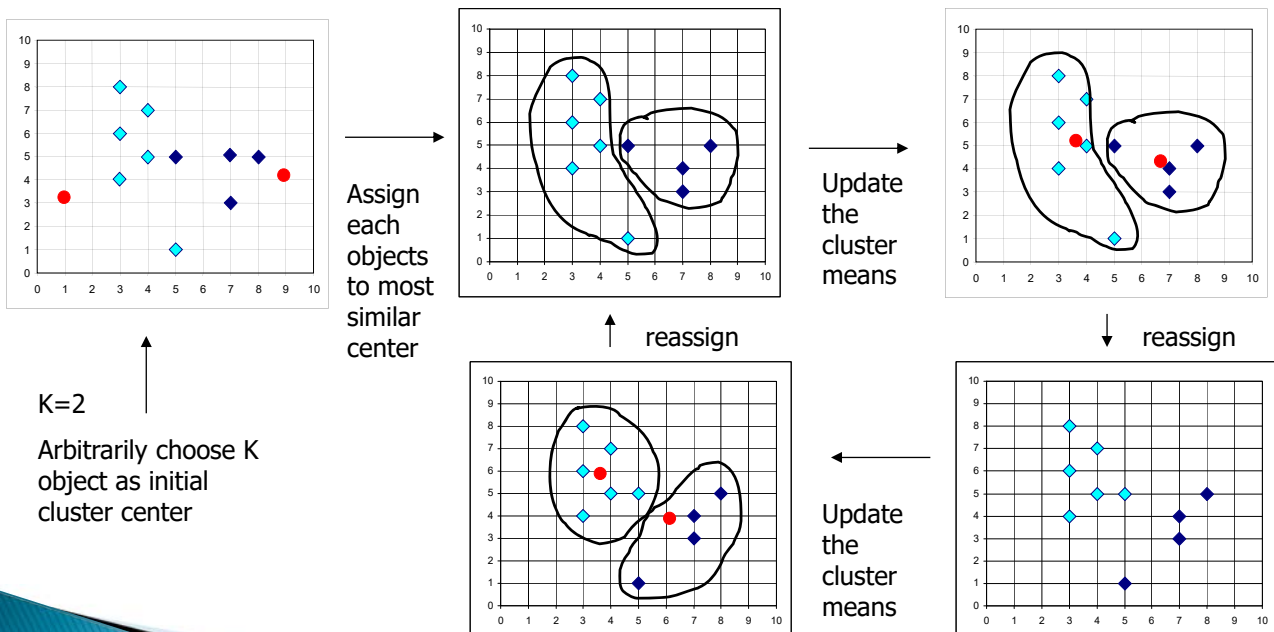
- ▶ Given a k , find a partition of k clusters that optimizes the chosen partitioning criterion
 - Global optimal: exhaustively enumerate all partitions
 - Heuristic methods: *k-means* and *k-medoids* algorithms
 - *k-means* (MacQueen'67): Each cluster is represented by the center of the cluster
 - *k-medoids* or PAM (Partition around medoids) (Kaufman & Rousseeuw'87): Each cluster is represented by one of the objects in the cluster

K-Means

- ▶ 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 partition (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 no more new assignment

K-Means

▶ Example



Comments

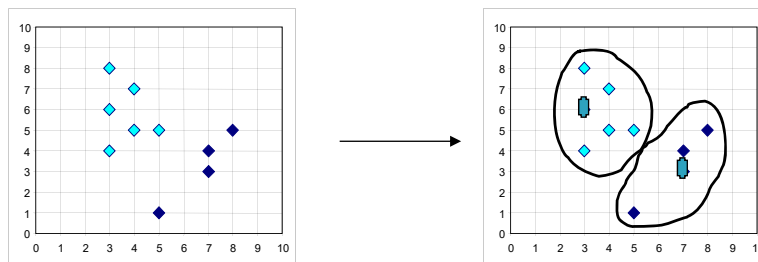
- ▶ Strength: *Relatively efficient*. $O(tkn)$, where n is # objects, k is # clusters, and t is # iterations. Normally, $k, t \ll n$.
 - Comparing: PAM: $O(k(n-k)^2)$, CLARA: $O(ks^2 + k(n-k))$
- ▶ Comment: Often terminates at a *local optimum*. The *global optimum* may be found using techniques such 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*

Variations of *K-Means*

- ▶ A few 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* (Huang'98)
 - Replacing means of clusters with modes
 - 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

Problem of K-Means

- ▶ 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.



K-Medoids

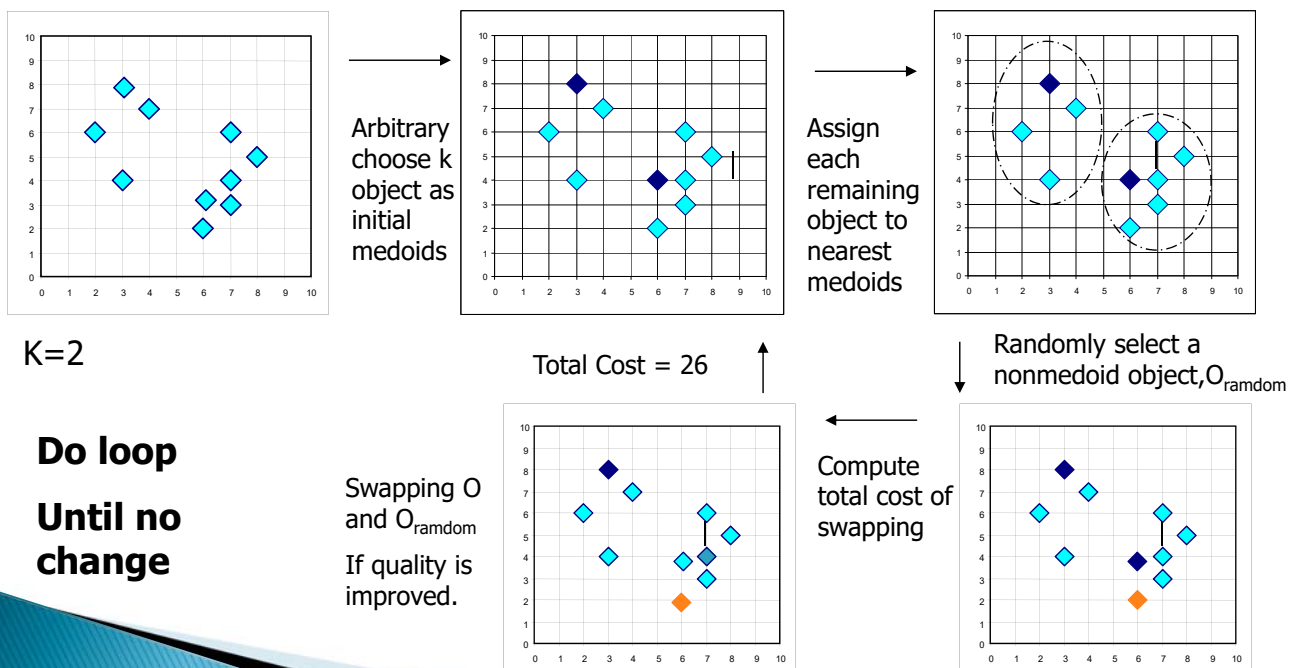
- ▶ Find *representative* objects, called medoids, in clusters
- ▶ *PAM* (Partitioning Around Medoids, 1987)
 - starts from an initial set of medoids and iteratively replaces one of the medoids by one of the non-medoids if it improves the total distance of the resulting clustering
 - *PAM* works effectively for small data sets, but does not scale well for large data sets
- ▶ *CLARA* (Kaufmann & Rousseeuw, 1990)
- ▶ *CLARANS* (Ng & Han, 1994): Randomized sampling
- ▶ Focusing + spatial data structure (Ester et al., 1995)

PAM (Partitioning Around Medoids)

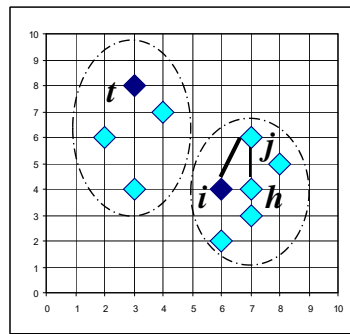
- ▶ PAM (Kaufman and Rousseeuw, 1987)
- ▶ Use real object to represent the cluster
 1. Select k representative objects arbitrarily
 2. For each pair of non-selected object h and selected object i , calculate the total swapping cost TC_{ih}
 3. For each pair of i and h ,
 - If $TC_{ih} < 0$, i is replaced by h
 - Then assign each non-selected object to the most similar representative object
 4. repeat steps 2–3 until there is no change

PAM

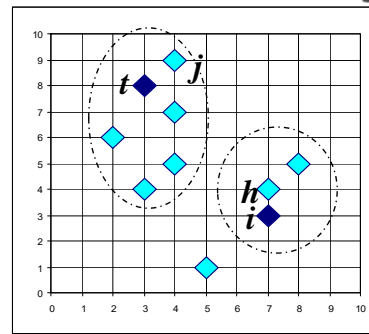
▶ Example



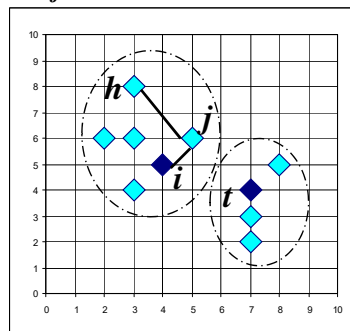
Total swapping cost $TC_{ih} = \sum_j C_{jih}$



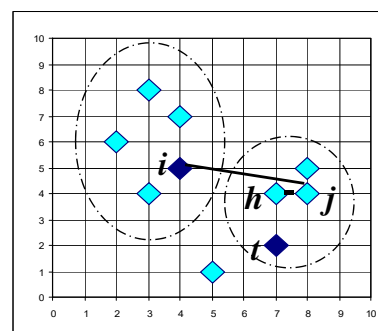
$$C_{jih} = d(j, h) - d(j, i)$$



$$C_{jih} = 0$$



$$C_{jih} = d(j, t) - d(j, i)$$



$$C_{jih} = d(j, h) - d(j, t)$$

Data Mining & Social Network Analysis 2021/03/24

35

Problem with PAM

- ▶ Pam is more robust than k-means in the presence of noise and outliers because a medoid is less influenced by outliers or other extreme values than a mean
- ▶ Pam works efficiently for small data sets but does not **scale well** for large data sets.
 - $O(k(n-k)^2)$ for each iteration

where n is # of data, k is # of clusters

➔ Sampling based method,

CLARA(Clustering LARge Applications)

CLARA (Clustering Large Applications)

- ▶ *CLARA* (Kaufmann and Rousseeuw in 1990)
- ▶ It draws *multiple samples* of the data set, applies *PAM* on each sample, and gives the best clustering as the output
- ▶ Strength: deals with larger data sets than *PAM*
- ▶ Weakness:
 - Efficiency depends on the sample size
 - A good clustering based on samples will not necessarily represent a good clustering of the whole data set if the sample is biased

CLARANS (“Randomized” CLARA)

- ▶ *CLARANS* (A Clustering Algorithm based on Randomized Search) (Ng and Han’94)
- ▶ *CLARANS* draws sample of neighbors dynamically
- ▶ The clustering process can be presented as searching a graph where every node is a potential solution, that is, a set of k medoids
- ▶ If the local optimum is found, *CLARANS* starts with new randomly selected node in search for a new local optimum
- ▶ It is more efficient and scalable than both *PAM* and *CLARA*
- ▶ Focusing techniques and spatial access structures may further improve its performance (Ester et al.’95)

References

- ▶ Slides from Prof. J.-W. Han, UIUC
- ▶ Slides from Prof. M.-S. Chen, NTU
- ▶ Slides from Prof. W.-Z. Peng, NCTU