



CLUSTERING

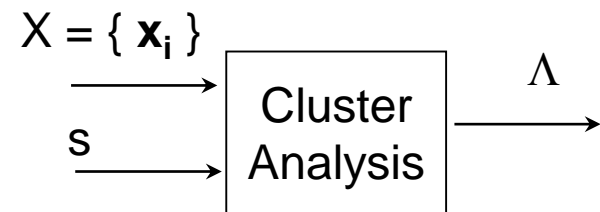
Bor-shen Lin

bslin@cs.ntust.edu.tw

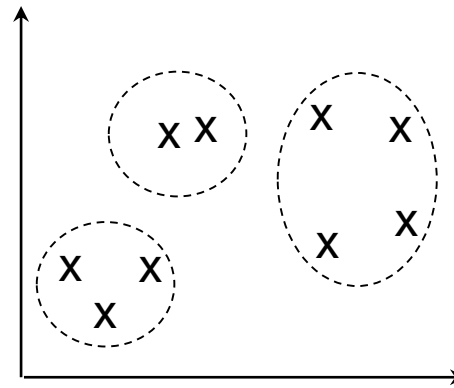
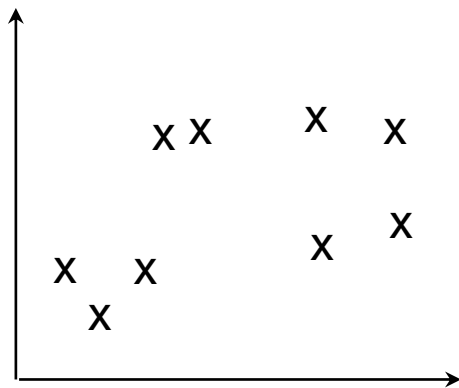
<http://www.cs.ntust.edu.tw/~bslin>

CLUSTER ANALYSIS

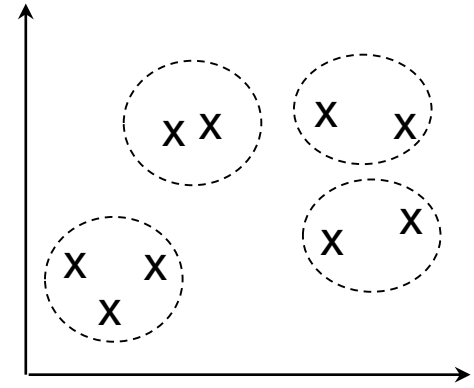
- A form of unsupervised learning
- Automatic classification of samples into a number of groups using a measure of association
 - Data in each group are **similar**
- Input
 - A set of samples: X
 - A measure of similarity: s
- Output
 - A number of groups
 - A partition $\Lambda = \{ G_1, G_2, \dots, G_N \}$
 $G_1 \cup G_2 \cup \dots \cup G_N = X$ and $G_i \cap G_j = \phi$



EXAMPLE



$K = 3$



$K = 4$

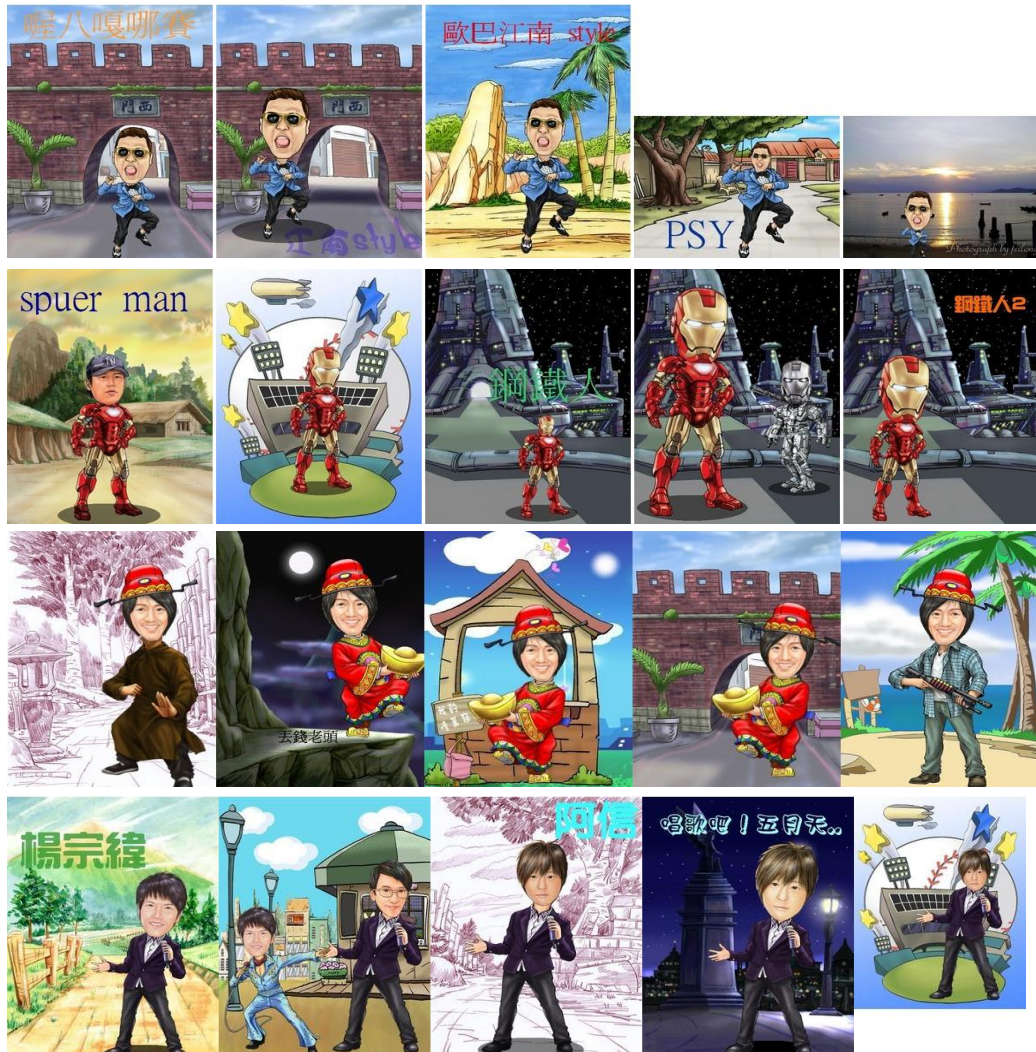


CLUSTERING FOR ANYTHING

- Clustering for customers
- Clustering for documents
- Clustering for words
- Clustering for images
- Clustering for melodies
- Clustering for objects
- Clustering for design patterns
- Clustering for styles



CLUSTERING OF DESIGN PATTERNS



BASIC CONCEPT OF CLUSTER ANALYSIS

- Sample
 - A point in a multi-dimensional space
- Objective of Cluster Analysis
 - Construct decision boundaries (classification surfaces) based on **unlabeled** training data set (Unsupervised learning)
 - Appropriate for exploration of interrelationships among samples to **make a preliminary assessment of the sample structure**



DIFFICULTIES OF CLUSTERING

- Shape and size of data
- Number of clusters
 - Determined according to desired resolution
 - There might be different results for the same data (different # of clusters)
- Human can visualize the data with dimension at most 3. Visualization of high dimensional data need to rely on clustering algorithms



FEATURES

- Quantitative features
 - Continuous values: \mathbb{R} , \mathbb{R}^+
 - Discrete values: $\{0, 1\}$, \mathbb{Z}
 - Interval values: $\{x \leq 20, 20 < x \leq 40, x > 40\}$
- Qualitative features
 - Nominal or unordered: color is “blue” or “red”
 - Ordinal: military rank: “general” and “colonel”



SIMILARITY/DISTANCE MEASURE

○ Similarity

- Symmetric: $s(\mathbf{x}, \mathbf{x}') = s(\mathbf{x}', \mathbf{x})$ for all \mathbf{x}, \mathbf{x}' in X
- Normalized: $0 \leq s(\mathbf{x}, \mathbf{x}') \leq 1$

○ Distance (dissimilarity) measure

- $d(\mathbf{x}, \mathbf{x}') \geq 0$ ($d(\mathbf{x}, \mathbf{x}) = 0$)
- $d(\mathbf{x}, \mathbf{x}') = d(\mathbf{x}', \mathbf{x})$
- Called **metric distance measure** if **triangular inequality** holds:
$$d(\mathbf{x}, \mathbf{x}'') \leq d(\mathbf{x}, \mathbf{x}') + d(\mathbf{x}', \mathbf{x}'')$$



METRICS FOR CONTINUOUS VALUES

- Euclidean distance

- $d_2(\mathbf{x}_i, \mathbf{x}_j) = [\sum_k (x_{ik} - x_{jk})^2]^{1/2}$

- Block distance (L1 metric)

- $d_1(\mathbf{x}_i, \mathbf{x}_j) = \sum_k |x_{ik} - x_{jk}|$

- Minkowski metric

- $d_p(\mathbf{x}_i, \mathbf{x}_j) = [\sum_k |x_{ik} - x_{jk}|^p]^{1/p}$

- Cosine similarity

- $s_{\cos}(\mathbf{x}_i, \mathbf{x}_j) = [\sum_k (x_{ik} \cdot x_{jk})] / [\sum_k x_{ik}^2 \cdot \sum_k x_{jk}^2]^{1/2}$

- $s_{\cos}(\mathbf{x}_i, \mathbf{x}_j) = 1$ if $\mathbf{x}_i = c \cdot \mathbf{x}_j, c > 0$



METRICS FOR DISCRETE VALUES

- $\mathbf{x}_i, \mathbf{x}_j$

$(1, 1) : a, (1, 0) : b$

$(0, 1) : c, (0, 0) : d$

- **Simple Matching Coefficient (SMC)**

- $S_{\text{smc}}(\mathbf{x}_i, \mathbf{x}_j) = (a + d) / (a + b + c + d)$

- **Jaccard Coefficient**

- $S_{\text{jc}}(\mathbf{x}_i, \mathbf{x}_j) = a / (a + b + c)$

- **Rao's Coefficient**

- $S_{\text{rc}}(\mathbf{x}_i, \mathbf{x}_j) = a / (a + b + c + d)$



EXAMPLE

- $\mathbf{x}_1 = (0, 0, 1, 1, 0, 1, 0, 1)$

$$\mathbf{x}_2 = (0, 1, 1, 0, 0, 1, 0, 0)$$

$$\rightarrow a = 2 \text{ for } (1, 1)$$

$$b = 2 \text{ for } (1, 0)$$

$$c = 1 \text{ for } (0, 1)$$

$$d = 3 \text{ for } (0, 0)$$

$$\rightarrow S_{\text{smc}}(\mathbf{x}_1, \mathbf{x}_2) = 5 / 8$$

$$S_{\text{jc}}(\mathbf{x}_1, \mathbf{x}_2) = 2 / 5$$

$$S_{\text{rc}}(\mathbf{x}_1, \mathbf{x}_2) = 2 / 8$$

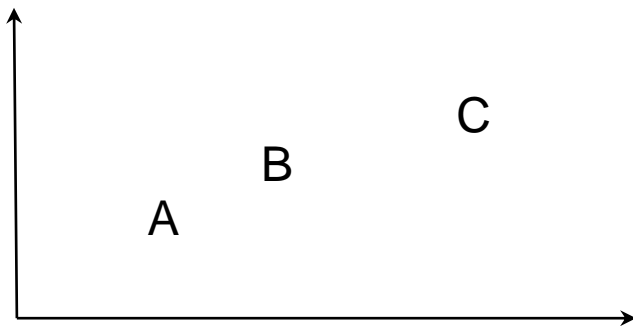


ADVANCED DISTANCE MEASURES

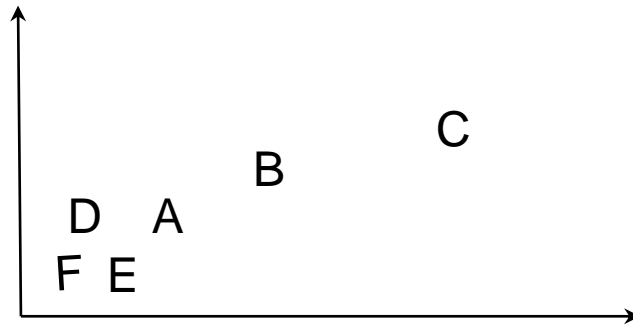
- Consider the effect of the context (relative distance)
- Mutual neighbor distance (MND)

$$\text{MND}(\mathbf{x}_i, \mathbf{x}_j) = \text{NN}(\mathbf{x}_i, \mathbf{x}_j) + \text{NN}(\mathbf{x}_j, \mathbf{x}_i)$$

$\text{NN}(\mathbf{x}_i, \mathbf{x}_j)$: the neighbor number of \mathbf{x}_i with respect to \mathbf{x}_j
 \mathbf{x}_i is the n -th closest point for \mathbf{x}_j



$\text{NN}(A, B) = 1$, $\text{NN}(B, A) = 1$
 $\text{NN}(B, C) = 1$, $\text{NN}(C, B) = 2$
 $\text{MND}(A, B) = 2$, $\text{MND}(B, C) = 3$



$\text{NN}(A, B) = 1$, $\text{NN}(B, A) = 4$
 $\text{NN}(B, C) = 1$, $\text{NN}(C, B) = 2$
 $\text{MND}(A, B) = 5$, $\text{MND}(B, C) = 3$

Positions of A & B are the same, but with different distances!



CLUSTERING ALGORITHMS

- Hierarchical Clustering Algorithm
 - Producing the hierarchy (dendrogram)
- Partitional Clustering algorithm
 - Producing the partition of the data



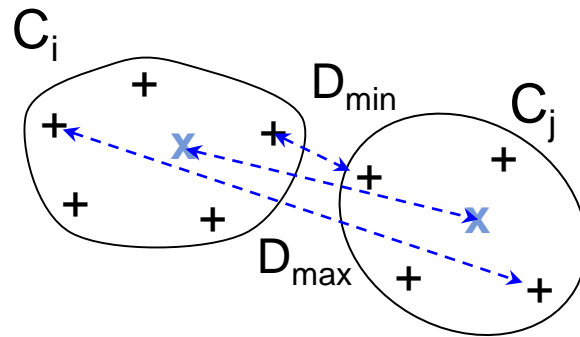
HIERARCHICAL CLUSTERING ALGORITHMS

- Divisible algorithm
 - A Top-down process
 - Example: CART for clustering
 - Divided by discrete properties
 - Starts from the entire set of samples
 - Divides it into a partition of subsets
- Agglomerative algorithm (aggregation)
 - A bottom-up process
 - Rain drops (small → big)
 - Regards each object as a cluster initially
 - The clusters are merged into larger clusters
 - A **dendrogram** is constructed



DISTANCE BETWEEN CLUSTERS

- $D_{\min}(C_i, C_j) = \min |\mathbf{x}_i - \mathbf{x}_j|$, for \mathbf{x}_i in C_i , \mathbf{x}_j in C_j
- $D_{\text{mean}}(C_i, C_j) = |\mathbf{m}_i - \mathbf{m}_j|$, for centroids \mathbf{m}_i and \mathbf{m}_j
- $D_{\text{avg}}(C_i, C_j) = 1/(n_i n_j) \sum_i \sum_j |\mathbf{x}_i - \mathbf{x}_j|$, for \mathbf{x}_i in C_i , \mathbf{x}_j in C_j
- $D_{\max}(C_i, C_j) = \max |\mathbf{x}_i - \mathbf{x}_j|$, for \mathbf{x}_i in C_i , \mathbf{x}_j in C_j



AGGLOMERATIVE CLUSTERING

1. Place each sample in its own cluster. Construct a list of **inter-cluster distances** for all pairs of samples, and sort this list in ascending order.
2. Step through the sorted list of distances, forming for each distinct threshold value d_k a graph of the samples where **pairs of samples closer than d_k are connected** into a new cluster by a graph edge. If all the samples are members of a connected graph, stop. Otherwise, repeat this step.
3. The output of the algorithm is a nested hierarchy of graphs, which **can be cut at the desired dissimilarity level** forming a partition (clusters) identified by simple connected components in the corresponding subgraph.



EXAMPLE

- Samples

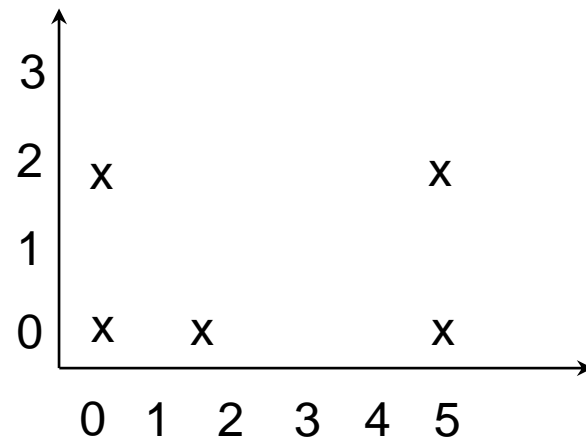
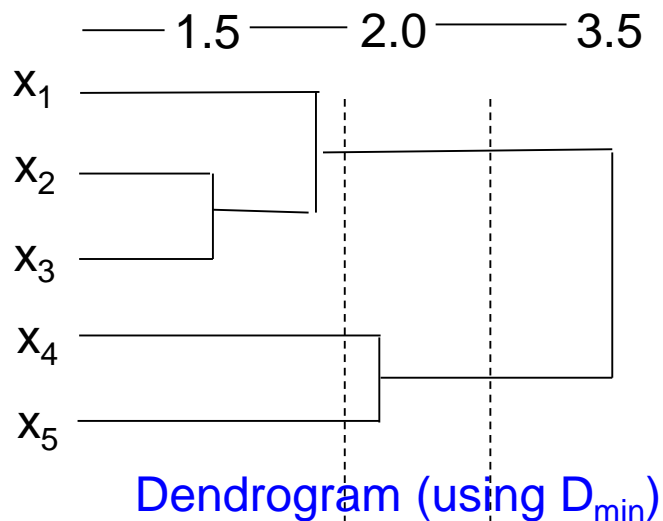
- $x_1=(0,2)$, $x_2=(0,0)$, $x_3=(1.5,0)$, $x_4=(5,0)$, $x_5=(5,2)$

- $d(x_1, x_2) = 2$, $d(x_1, x_3) = 2.5$, $d(x_1, x_4) = 5.39$,

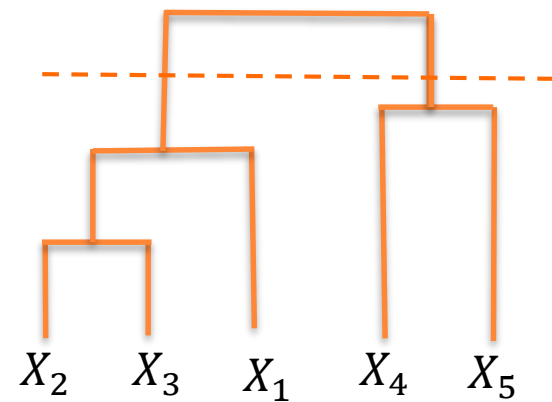
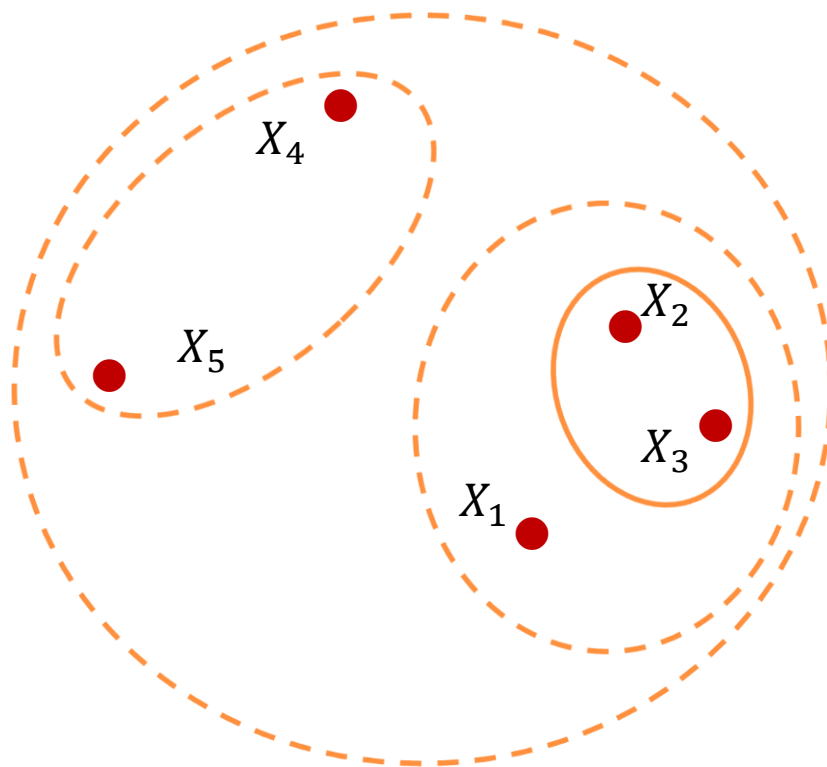
- $d(x_1, x_5) = 5$, $d(x_2, x_3) = 1.5$, $d(x_2, x_4) = 5$,

- $d(x_2, x_5) = 5.29$, $d(x_3, x_4) = 3.5$,

- $d(x_3, x_5) = 4.03$, $d(x_4, x_5) = 2$



MERGING OF CLUSTERS



CHAMELEON ALGORITHM

- Basic idea: Improve the clustering quality by using a more elaborate criterion when merging two clusters
 - C_i and C_j will be merged if the **interconnectivity and closeness** of the merged cluster is very similar to that of C_i and C_j before merging
- Min-cut on a graph (of cluster C_i)
 - Partitioning a graph into **two** parts of close, equal size such that **the total weight of the edges being cut** is minimized.
 - Total weight of the edges: stand for total interconnectivity (c.f. inner links in political parties)
 - min-cut: minimize loss of connectivity for the graph after being cut
 - Interconnectivity $I(C_i)$: **total weight of edges being cut (total loss)**
 - Closeness $C(C_i)$: **average weight of edges being cut (average loss)**



CHAMELEON ALGORITHM

1. Divide all data into *highly dense* sub-clusters.
 - Construct initial graph $G=(V,E)$ by KNN with weighted edge $e(v_i,v_j)$: closeness between two samples
 - Recursively partition G into many *small, unconnected subgraphs* by doing min-cut
 - Stopped when certain criteria are satisfied
2. Merge sub-clusters into larger clusters.
 - For C_i and C_j , compute $RI(C_i, C_j)$ and $RC(C_i, C_j)$
 - Relative interconnectivity
$$RI(C_i, C_j) \equiv I(C_i \cup C_j) / [0.5 \cdot (I(C_i) + I(C_j))]$$
 - Relative closeness
$$RC(C_i, C_j) \equiv C(C_i \cup C_j) / [0.5 \cdot (C(C_i) + C(C_j))]$$
 - $s(C_i, C_j) \equiv RC(C_i, C_j) \cdot RI(C_i, C_j)^\alpha, 0 \leq \alpha \leq 1$
Minimizing the decrease after being merged



PARTITIONAL CLUSTERING

- Why?
 - Construction of dendrogram is computationally complex for large data set
- Global criterion
 - Euclidean square-error measure E^2 for k-th cluster
 $E^2 = \sum_k d_k$, $d_k = \sum_i (x_{ik} - M_k)^2$, $\mathbf{m}_k = (1/n_k) \sum_i \mathbf{x}_{ik}$
 - Euclidean distance cannot be used if the dimensions have different physical meaning → using stochastic model
- Local criterion
 - Minimal **mutual neighbor distance** (MND)
 - Utilizing the local structure or context in the data

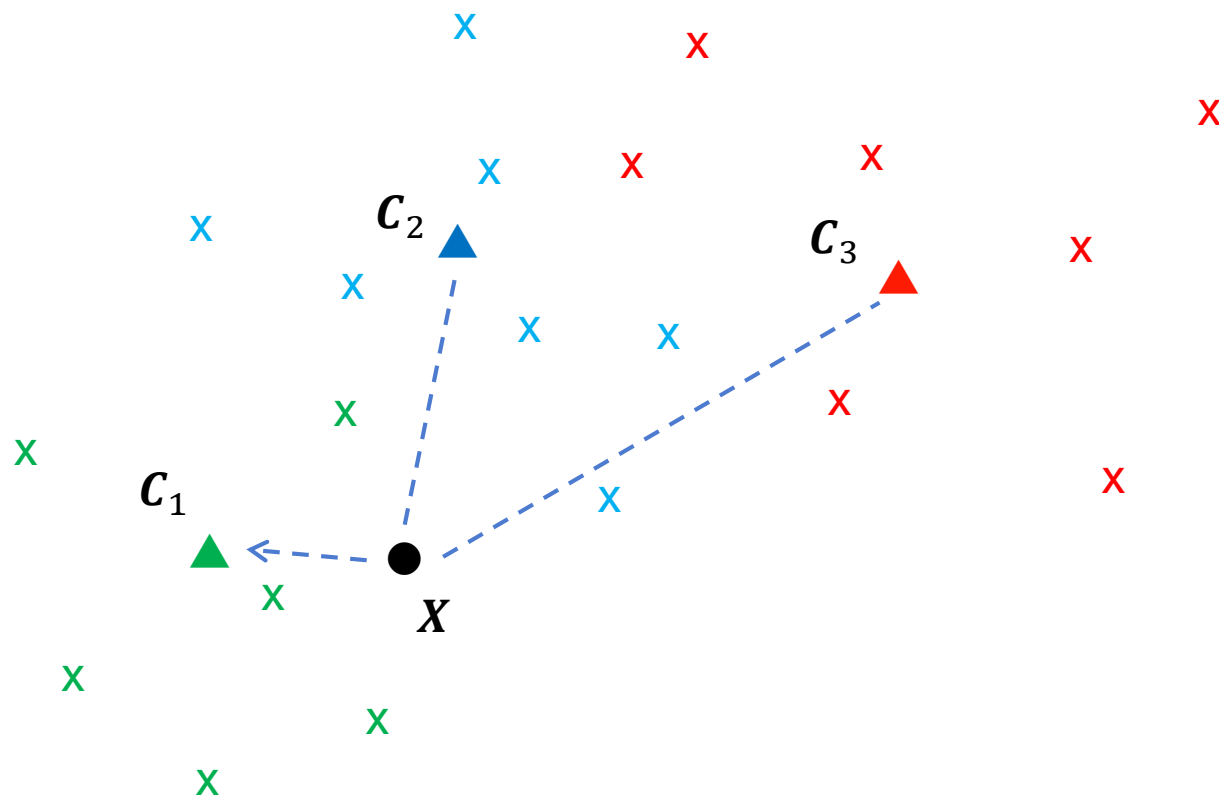


K-MEANS CLUSTERING ALGORITHM

- A partitional clustering algorithm employing a square-error criterion
- Simplest & most commonly used
 - Easy to implement
 - Time and space complexity is relatively small
 - Gives surprisingly good result if the clusters are compact, hyperspherical in shape, and well separated in the feature space
- Work well for data that are far apart
- Equivalent to Kohonen net in ANN
- Similar methods (e.g. C-means, K-mediods)



REDISTRIBUTION OF SAMPLES



K-MEANS ALGORITHM

1. Select an *initial partition* with k clusters containing randomly chosen samples, and compute the centroids of the clusters.
2. Generate a new partition by (re-distributing) assigning each sample to its closest cluster center $k^* = \operatorname{argmin}_k d(\mathbf{x}, \mathbf{m}_k)$ for every \mathbf{x}
3. Compute new centroids of the clusters $\mathbf{m}_k = (1/n) \sum \mathbf{x}_i$ for cluster k
4. Repeat steps 2 and 3 until an optimum value of the criterion function is found.



CHARACTERISTICS OF K-MEANS

- Time complexity: $O(nkl)$
 - n: number of samples
 - k: number of clusters
 - l: number of iterations
- Space complexity: $O(k+n)$
- Order-independent
 - It generates the same partition of the data at the end of the partitioning process **irrespective of the order in which the samples are presented to the algorithm**



ORDER-DEPENDENT INCREMENTAL CLUSTERING

1. Initially, there is no cluster.
 2. For every point, compute the distance between the point and any existing cluster.
 3. If the minimum distance does exist, raise a new cluster with the point and go to step 2. Otherwise, go to step 3.
 4. If the minimum distance is smaller than a threshold, add that point to the corresponding cluster. Otherwise raise a new cluster with the point.
 5. Go to step 2.
- Order dependent
 - The order that data are presented influences the result



MAJOR PROBLEMS OF K-MEANS

- Sensitive to the selection of the initial partition
- May converge to a local minimum
- Lack of available guidelines for
 - Choosing the initial partition
 - Adjusting the number of clusters
 - Selecting the stopping criterion
- Very sensitive to noise and outlier data points
 - Every point has the same contribution to the centroid
 - **K-medoids** method: remove the outliers when computing the centroids
 - less sensitive to noise and outliers
- Based on global distance
 - Might not work well for the data with locally connected patterns



K-MEANS FOR CATEGORICAL DATA

- Centroids for clusters can NOT be computed
 - $d(\mathbf{x}, \mathbf{m}_k)$ cannot be calculated without \mathbf{m}_k ,
but $d(\mathbf{x}_i, \mathbf{x}_j)$ or $s(\mathbf{x}_i, \mathbf{x}_j)$ can still be found
- Use *K-nearest neighbor (KNN)* for reclustering data
 - \mathbf{x} can be clustered according to the clusters of its **KNNs** through majority voting
 - Do not compute the distance between \mathbf{x} and centroid
 - My choice → according to the choices of K-best friends
- Use *the distribution of cluster* for reclustering data
 - $P(\mathbf{x} | C_k)$ could be trained for each cluster
 - Example: customer \mathbf{x} = [gender=male, age=20, occupation=engineer, income=900k, ...]



REDISTRIBUTION OF CATEGORICAL DATA (KNN)

- $\mathbf{x}_1 = (A, B, A, B, C, B)$
 $\mathbf{x}_2 = (A, A, A, B, A, B)$
 $\mathbf{x}_3 = (B, B, A, B, A, B)$
 $\mathbf{x}_4 = (B, C, A, B, B, A)$
 $\mathbf{x}_5 = (B, A, B, A, C, A)$
 $\mathbf{x}_6 = (A, C, B, A, B, B)$
- $C_1 = \{ \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3 \}, C_2 = \{ \mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_6 \}$
- For $\mathbf{y} = \{ A, C, A, B, C, A \}$
 $s_{\text{smc}}(\mathbf{y}, \mathbf{x}_1)=0.66, s_{\text{smc}}(\mathbf{y}, \mathbf{x}_2)=0.50, s_{\text{smc}}(\mathbf{y}, \mathbf{x}_3)=0.33$
 $s_{\text{smc}}(\mathbf{y}, \mathbf{x}_4)=0.66, s_{\text{smc}}(\mathbf{y}, \mathbf{x}_5)=0.33, s_{\text{smc}}(\mathbf{y}, \mathbf{x}_6)=0.33$
- KNN for $K=3$ are $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_4 \rightarrow$ voting C_1 for \mathbf{y}



K-MEANS VS. GMM TRAINING

○ K-Means

- Distance: Euclidean distance (for homogenous data)
- Objective function: square error
- Clusters have *mutually exclusive* data
- Each point has equal contribution (weight) to mean

$$\mathbf{m}_k = \frac{1}{n_k} \sum_{\mathbf{x}_i \in C_k} \mathbf{x}_i$$

○ GMM

- Similarity: Gaussian probability (distance is normalized by covariance $(\mathbf{x} - \boldsymbol{\mu}_k)^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)$)
- Objective function: expectation of log probabilities
- All data points contribute to each cluster

$$\boldsymbol{\mu}_k = \frac{\sum_i l_i(k) \mathbf{x}_i}{\sum_{i,k} l_i(k)}, \quad \boldsymbol{\Sigma}_k = \frac{\sum_i l_i(k) (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^t}{\sum_{i,k} l_i(k)}$$



K-MEANS FOR VECTOR QUANTIZATION

- Select a set of prototype vectors to represent a large vector space
- Vectors in vector space are encoded into a discrete symbols
- *Codebook*: a set of *codewords* (vectors)
- Example: data compression



EVALUATION FOR CLUSTER ANALYSIS

- Assessment of the data domain
 - estimating the data in advance(how many clusters, distance between clusters, ...)
- Cluster validity (verification on the results)
 - External assessment of validity
 - Compare the discovered structure to an a priori structure
 - Internal examination of validity
 - Determine if the discovered structure is intrinsically appropriate or meaningful
 - Relative test
 - Adjust the algorithms and parameters to produce different clusters. Compare the results to judge which is more applicable



SUMMARY

- There is NO clustering technique that is universally applicable in uncovering the variety of structures present in multidimensional data sets.
- User's understanding of the problem and the corresponding data types will be the best criteria to select the appropriate method.
- The data should be subjected to tests for clustering tendency before applying a clustering algorithm, followed by a validation of the clusters generated by the algorithm
- There is no best clustering algorithm. Try several ones.
- Clustering analysis might be used as processing techniques in intermediate level.
 - e.g. word class n-gram (sharing probabilities)



REFERENCE

- Data Mining: Concepts, Models, Methods and Algorithms

Mehmed Kantardzic, Wiley Inter-science

