

Data Mining and Data Warehousing

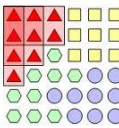
Chapter 5

Cluster Analysis

Instructor: Suresh Pokharel

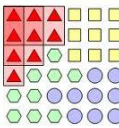
ME in ICT (Asian Institute of Technology, Thailand)

BE in Computer (NCIT, Pokhara University)



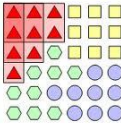
Cluster Analysis

- What is Cluster Analysis?
- Types of Data in Cluster Analysis
- A Categorization of Major Clustering Methods
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Grid-Based Methods
- Model-Based Clustering Methods
- Summary



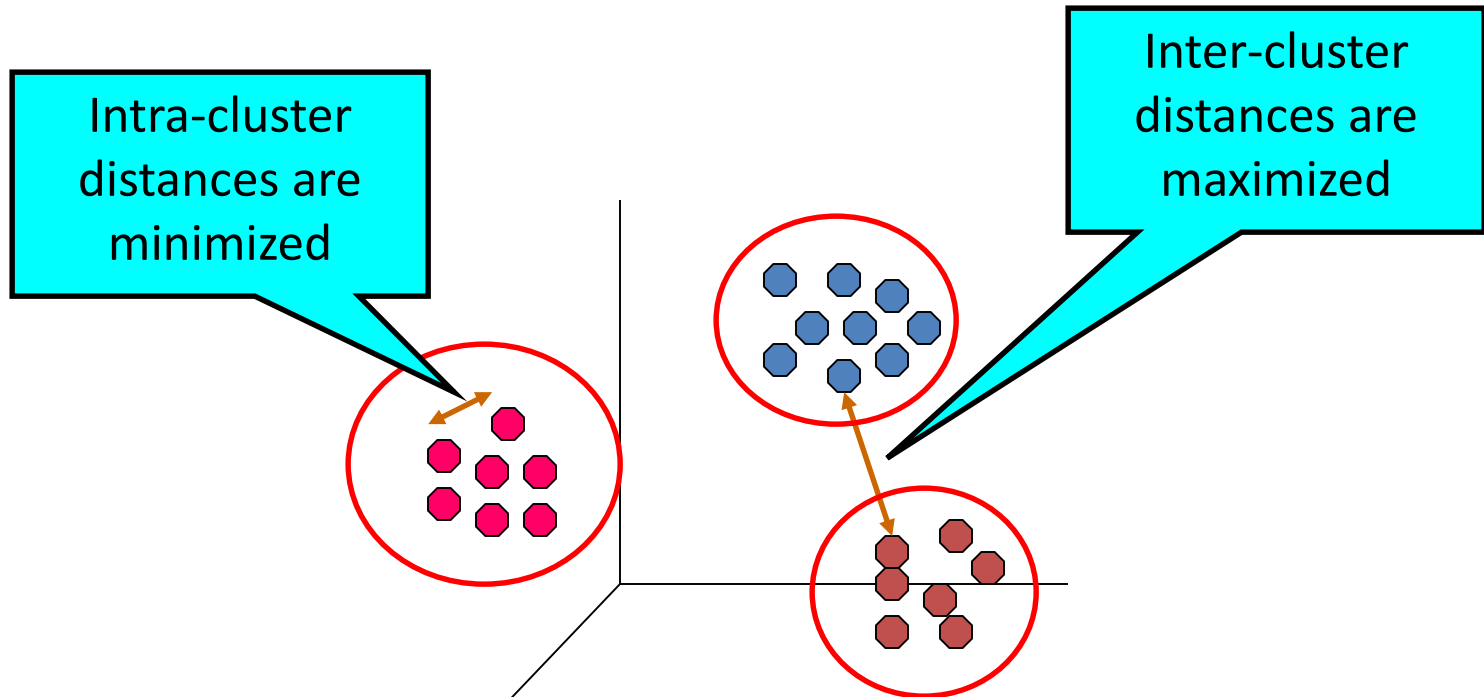
What is Cluster Analysis?

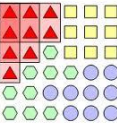
- Cluster: a collection of data objects
 - Similar to one another within the same cluster
 - Dissimilar to the objects in other clusters
- Cluster analysis
 - Grouping a set of data objects into clusters
- Clustering is **unsupervised classification**: no predefined classes
- Clustering is used:
 - As a **stand-alone tool** to get insight into data distribution
 - Visualization of clusters may unveil important information
 - As a **preprocessing step** for other algorithms
 - Efficient indexing or compression often relies on clustering



What is Cluster Analysis?

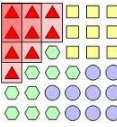
- Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups





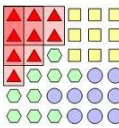
General Applications of Clustering

- Pattern Recognition
- Spatial Data Analysis
 - create thematic maps in GIS by clustering feature spaces
 - detect spatial clusters and explain them in spatial data mining
- Image Processing
 - cluster images based on their visual content
- Economic Science (especially market research)
- WWW and IR
 - document classification
 - cluster Weblog data to discover groups of similar access patterns



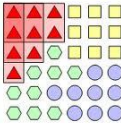
Examples of Clustering Applications

- 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



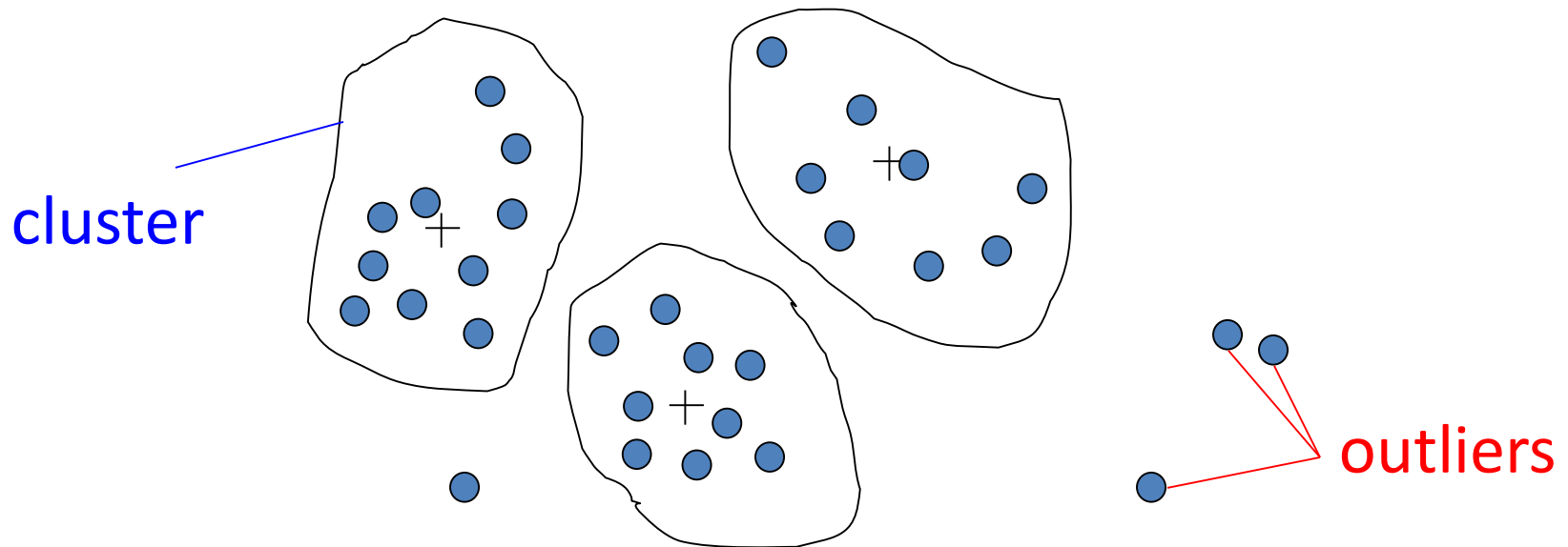
What Is Good Clustering?

- 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.



Outliers

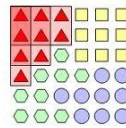
- Outliers are objects that do not belong to any cluster or form clusters of very small cardinality



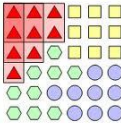
- In some applications we are interested in discovering outliers, not clusters (**outlier analysis**)



Cluster Analysis



- What is Cluster Analysis?
- Types of Data in Cluster Analysis
- A Categorization of Major Clustering Methods
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Grid-Based Methods
- Model-Based Clustering Methods
- Summary



Data Structures

- *data* matrix
– (two modes)

the “classic” data input

attributes/dimensions

tuples/objects

x_{11}	...	x_{1f}	...	x_{1p}
...
x_{i1}	...	x_{if}	...	x_{ip}
...
x_{n1}	...	x_{nf}	...	x_{np}

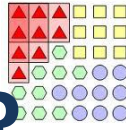
- *dissimilarity* or *distance*
matrix

the desired data input to some
clustering algorithms

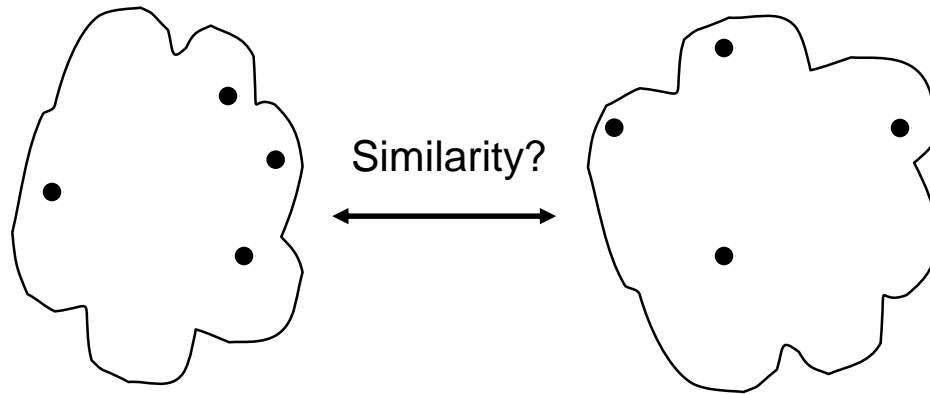
objects

objects

0			
$d(2,1)$	0		
$d(3,1)$	$d(3,2)$	0	
:	:	:	
$d(n,1)$	$d(n,2)$ 0



How to Define Inter-Cluster Similarity?



MIN

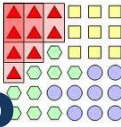
MAX

Group Average

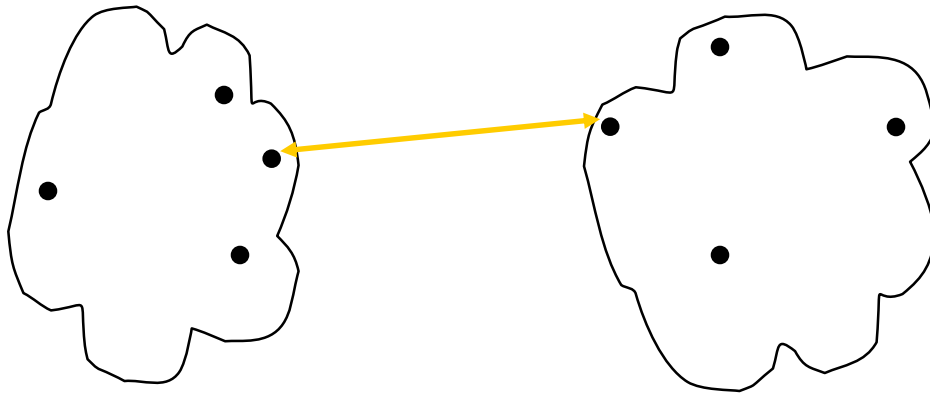
Distance Between Centroids

Other methods driven by an objective function

Ward's Method uses squared error



How to Define Inter-Cluster Similarity?



MIN

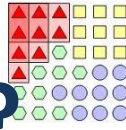
MAX

Group Average

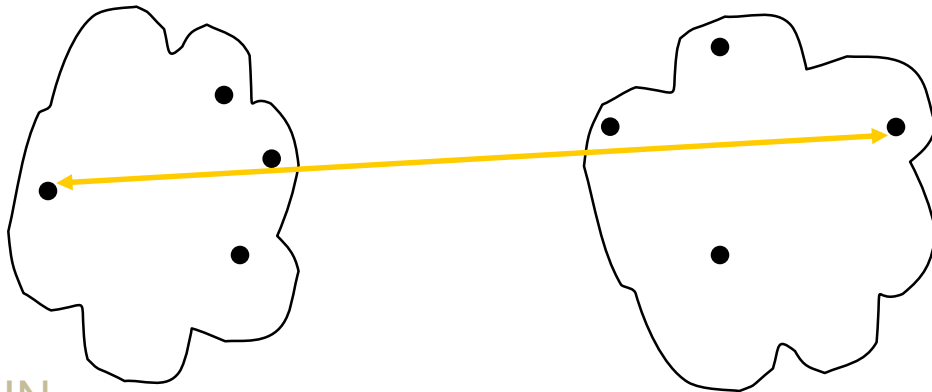
Distance Between Centroids

Other methods driven by an objective function

Ward's Method uses squared error



How to Define Inter-Cluster Similarity?



MIN

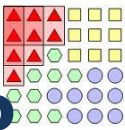
MAX

Group Average

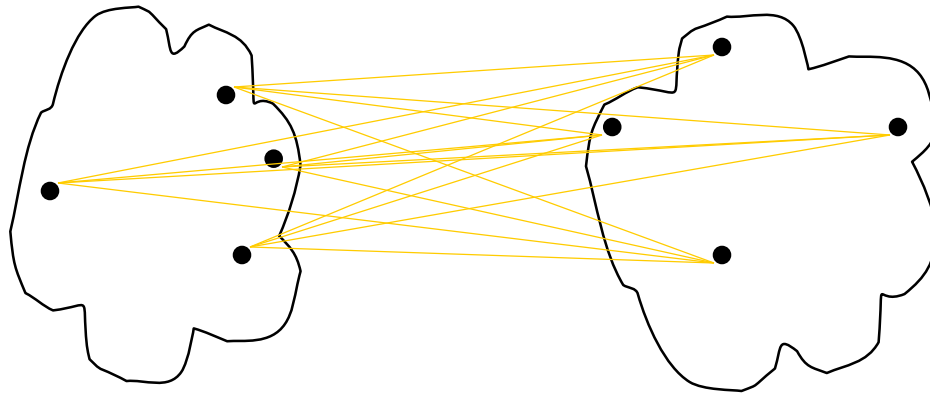
Distance Between Centroids

Other methods driven by an objective function

Ward's Method uses squared error



How to Define Inter-Cluster Similarity?



MIN

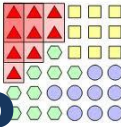
MAX

Group Average

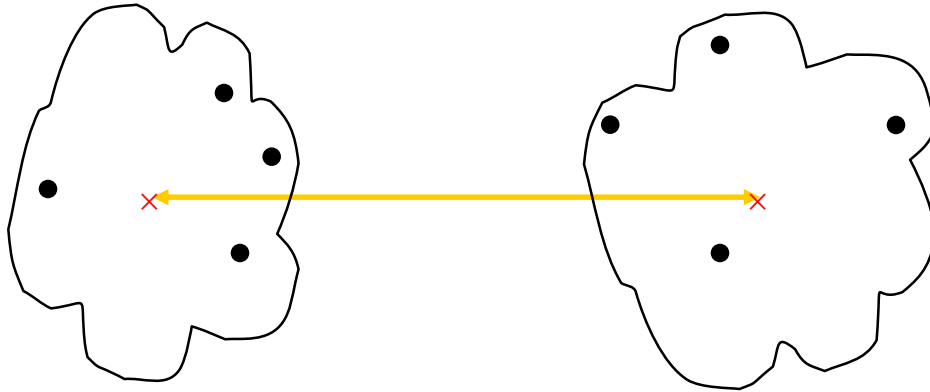
Distance Between Centroids

Other methods driven by an objective function

Ward's Method uses squared error



How to Define Inter-Cluster Similarity?



MIN

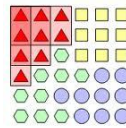
MAX

Group Average

Distance Between Centroids

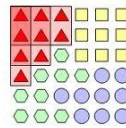
Other methods driven by an objective function

Ward's Method uses squared error



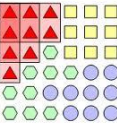
Cluster Analysis

- What is Cluster Analysis?
- Types of Data in Cluster Analysis
- A Categorization of Major Clustering Methods
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Grid-Based Methods
- Model-Based Clustering Methods
- Outlier Analysis
- Summary



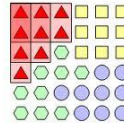
Major Clustering Approaches

- Partitioning algorithms: Construct random partitions and then iteratively refine them by some criterion
- Hierarchical algorithms: Create a hierarchical decomposition of the set of data (or objects) using some criterion
- Density-based: based on connectivity and density functions
- Grid-based: based on a multiple-level granularity structure
- Model-based: A model is hypothesized for each of the clusters and the idea is to find the best fit of that model to each other



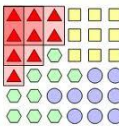
Cluster Analysis

- What is Cluster Analysis?
- Types of Data in Cluster Analysis
- A Categorization of Major Clustering Methods
- **Partitioning Methods**
- Hierarchical Methods
- Density-Based Methods
- Grid-Based Methods
- Model-Based Clustering Methods
- Summary



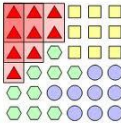
Partitioning Algorithms: Basic Concepts

- Partitioning method: Construct a partition of a database **D** of **n** objects into a set of **k** clusters
- 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



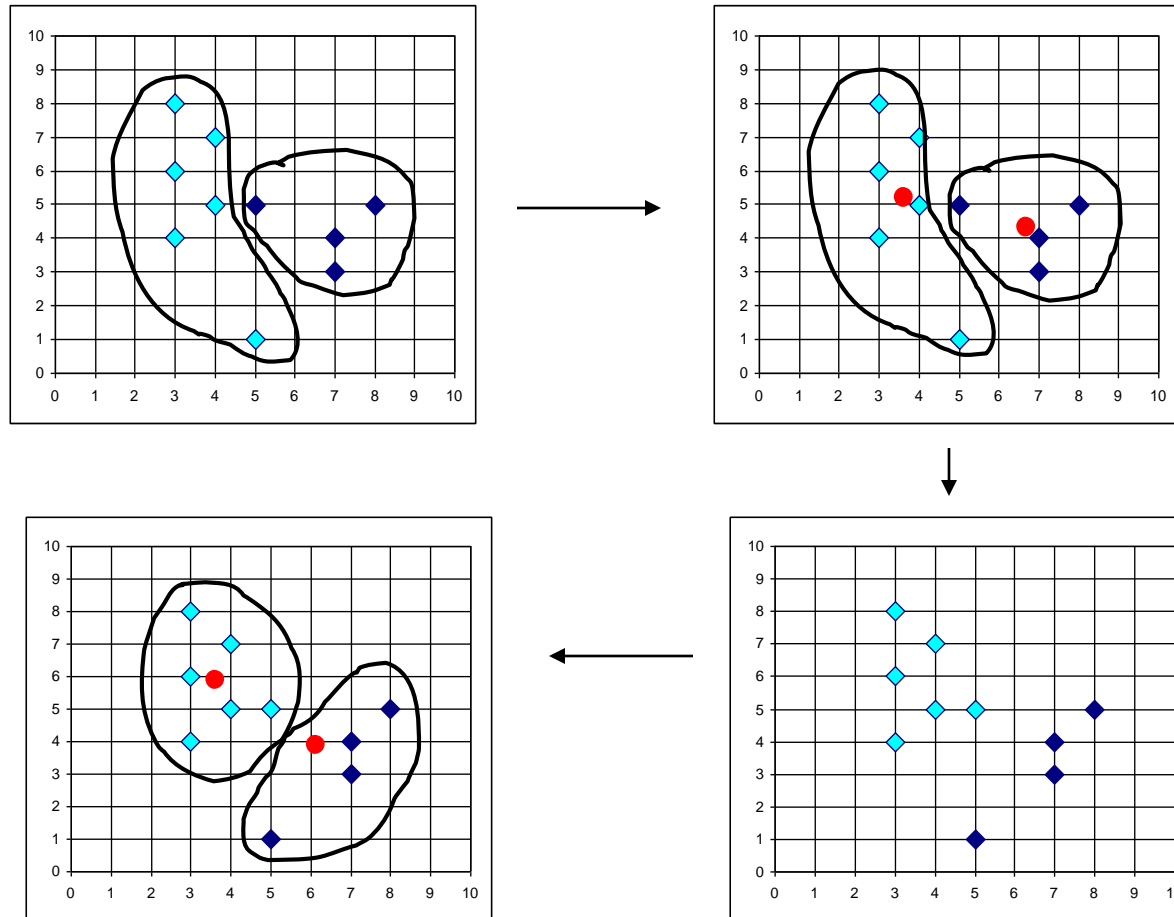
The k-means Clustering Method

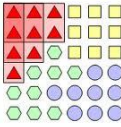
- Given k , the *k-means* algorithm is implemented in 4 steps:
 1. Partition objects into k nonempty subsets
 2. Compute seed points as the **centroids** of the clusters of the current partition. The centroid is the center (mean point) of the cluster.
 3. Assign each object to the cluster with the nearest seed point.
 4. Go back to Step 2, stop when no more new assignment.



The k-means Clustering Method

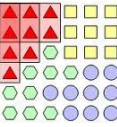
- Example





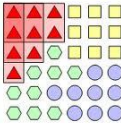
K-Means example

- 2, 3, 6, 8, 9, 12, 15, 18, 22 – break into 3 clusters
 - Cluster 1 - 2, 8, 15 – mean = 8.3
 - Cluster 2 - 3, 9, 18 – mean = 10
 - Cluster 3 - 6, 12, 22 – mean = 13.3
- Re-assign
 - Cluster 1 - 2, 3, 6, 8, 9 – mean = 5.6
 - Cluster 2 – mean = 0
 - Cluster 3 – 12, 15, 18, 22 – mean = 16.75
- Re-assign
 - Cluster 1 – 3, 6, 8, 9 – mean = 6.5
 - Cluster 2 – 2 – mean = 2
 - Cluster 3 = 12, 15, 18, 22 – mean = 16.75



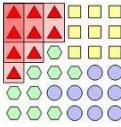
K-Means example (continued)

- Re-assign
 - Cluster 1 - 6, 8, 9 – mean = 7.6
 - Cluster 2 – 2, 3 – mean = 2.5
 - Cluster 3 – 12, 15, 18, 22 – mean = 16.75
- Re-assign
 - Cluster 1 - 6, 8, 9 – mean = 7.6
 - Cluster 2 – 2, 3 - mean = 2.5
 - Cluster 3 – 12, 15, 18, 22 – mean = 16.75
- No change, so we're done



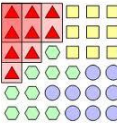
K-Means example – different starting order

- 2, 3, 6, 8, 9, 12, 15, 18, 22 – break into 3 clusters
 - Cluster 1 - 2, 12, 18 – mean = 10.6
 - Cluster 2 - 6, 9, 22 – mean = 12.3
 - Cluster 3 – 3, 8, 15 – mean = 8.6
- Re-assign
 - Cluster 1 - mean = 0
 - Cluster 2 – 12, 15, 18, 22 - mean = 16.75
 - Cluster 3 – 2, 3, 6, 8, 9 – mean = 5.6
- Re-assign
 - Cluster 1 – 2 – mean = 2
 - Cluster 2 – 12, 15, 18, 22 – mean = 16.75
 - Cluster 3 = 3, 6, 8, 9 – mean = 6.5



K-Means example (continued)

- Re-assign
 - Cluster 1 – 2, 3 – mean = 2.5
 - Cluster 2 – 12, 15, 18, 22 – mean = 16.75
 - Cluster 3 – 6, 8, 9 – mean = 7.6
- Re-assign
 - Cluster 1 – 2, 3 – mean = 2.5
 - Cluster 2 – 12, 15, 18, 22 - mean = 16.75
 - Cluster 3 – 6, 8, 9 – mean = 7.6
- No change, so we're done



Comments on the k-means Method

- Strength

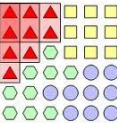
- *Relatively efficient: $O(tkn)$, where n is # objects, k is # clusters, and t is # iterations. Normally, $k, t \ll n$.*

- Weaknesses

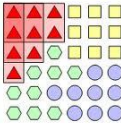
- Applicable only when *mean* is defined, then what about categorical data?
- Need to specify k , the *number* of clusters, in advance
- Unable to handle noisy data and *outliers*



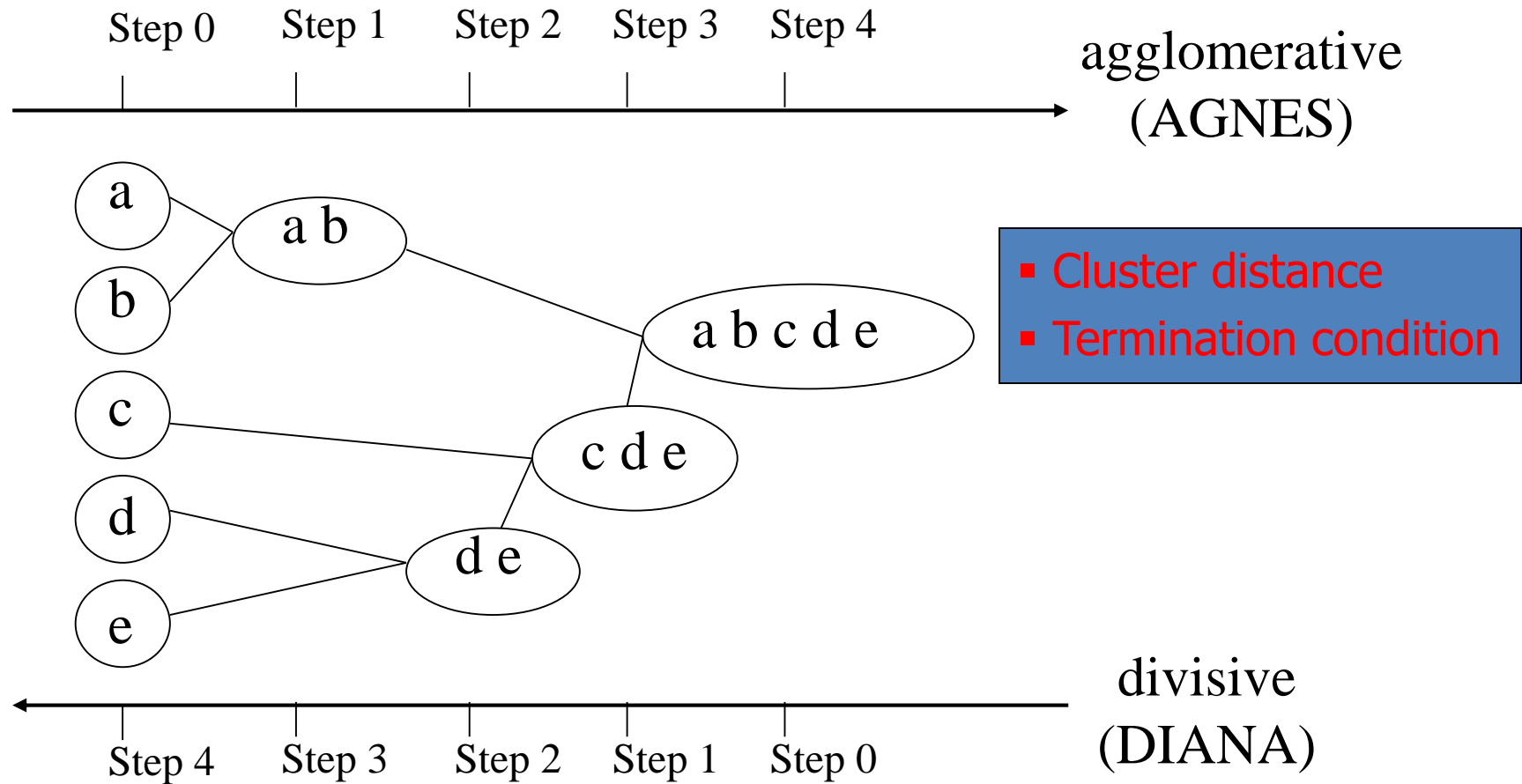
Cluster Analysis

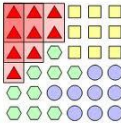


- What is Cluster Analysis?
- Types of Data in Cluster Analysis
- A Categorization of Major Clustering Methods
- Partitioning Methods
- **Hierarchical Methods**
- Density-Based Methods
- Grid-Based Methods
- Model-Based Clustering Methods
- Outlier Analysis
- Summary



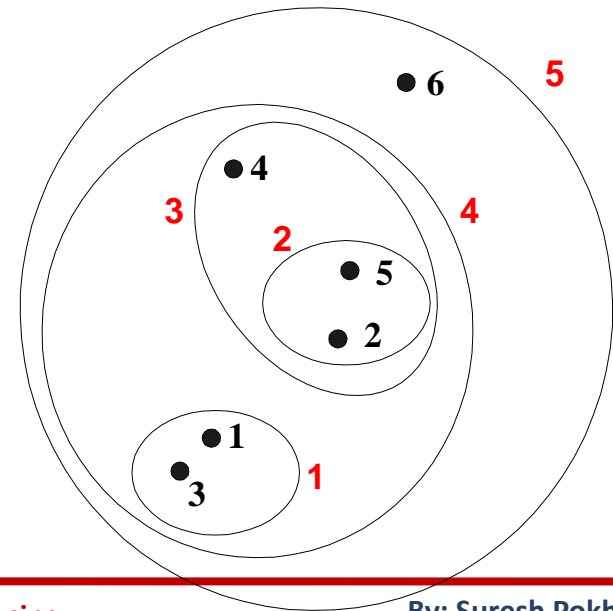
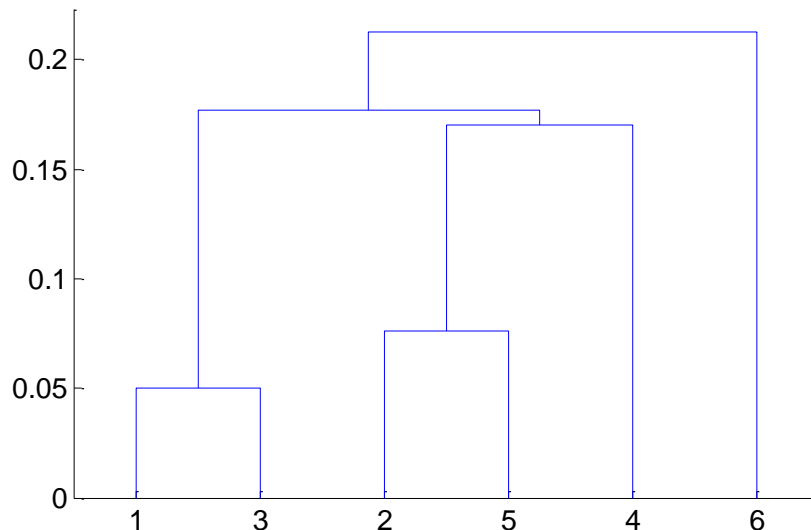
Hierarchical Clustering

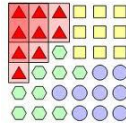




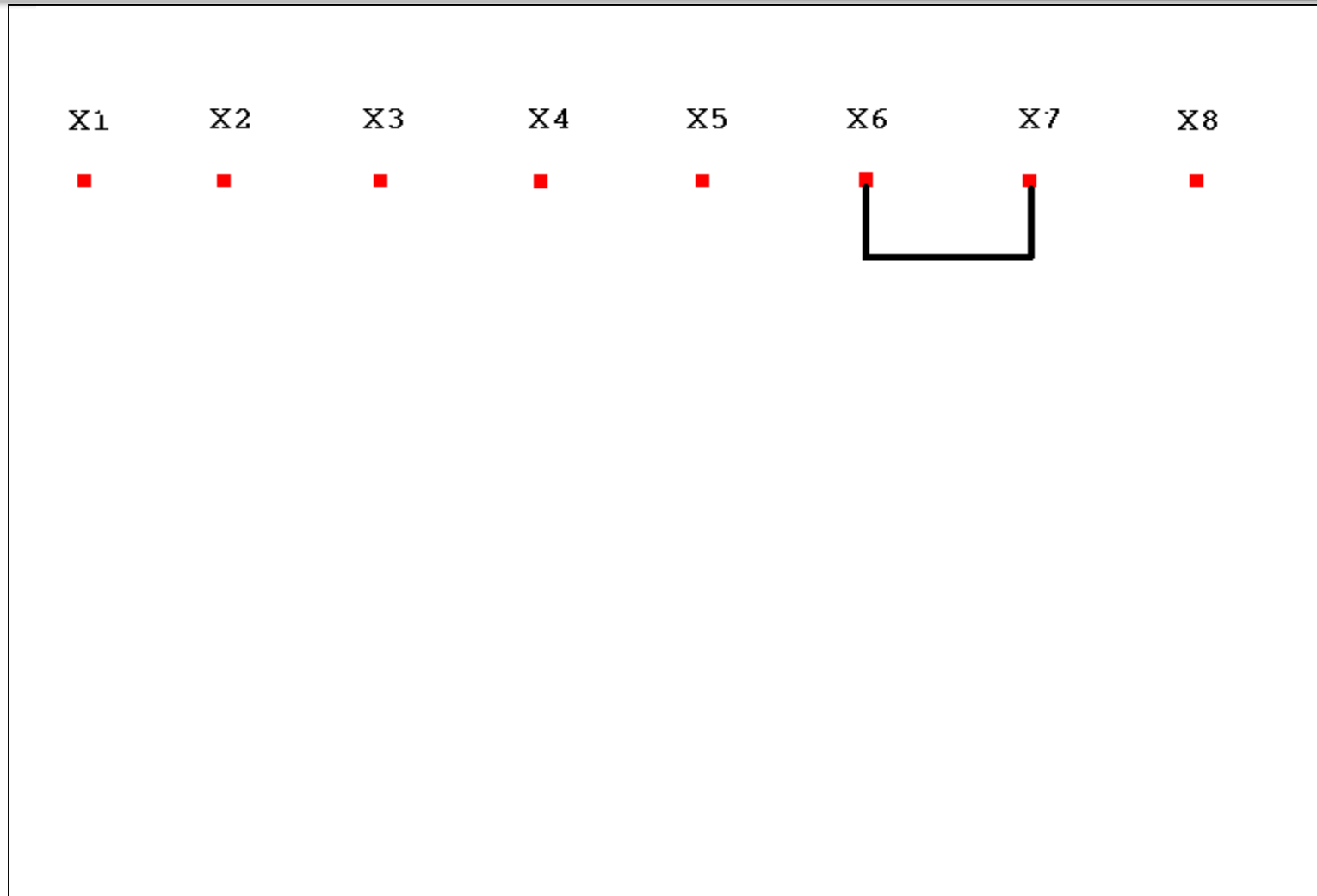
Hierarchical Clustering

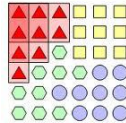
- Produces a set of nested clusters organized as a hierarchical tree
- Can be visualized as a “dendrogram”
 - A tree like diagram that records the sequences of merges or splits



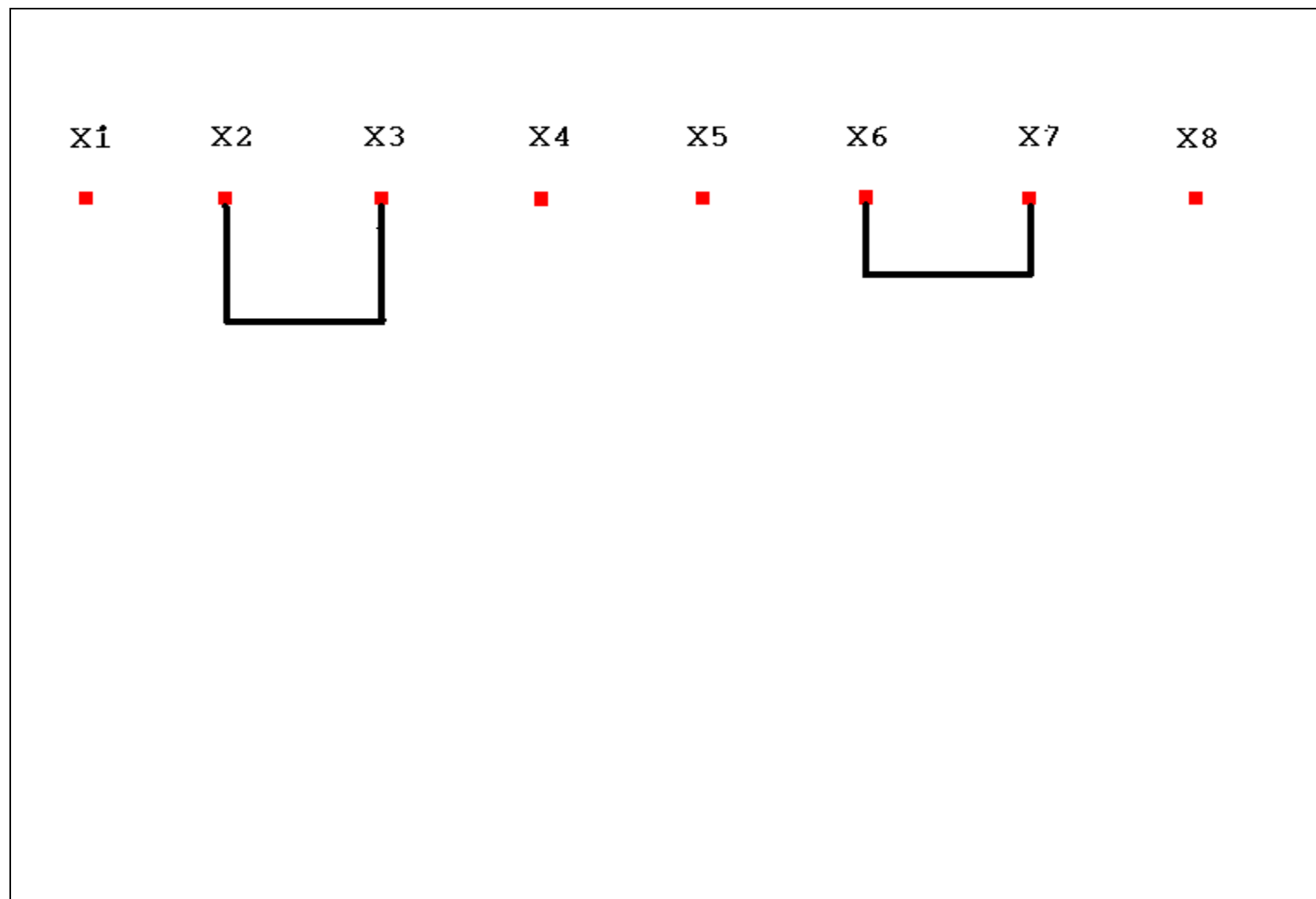


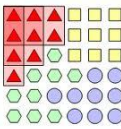
Nearest Neighbor, Level 2, $k = 7$ clusters.



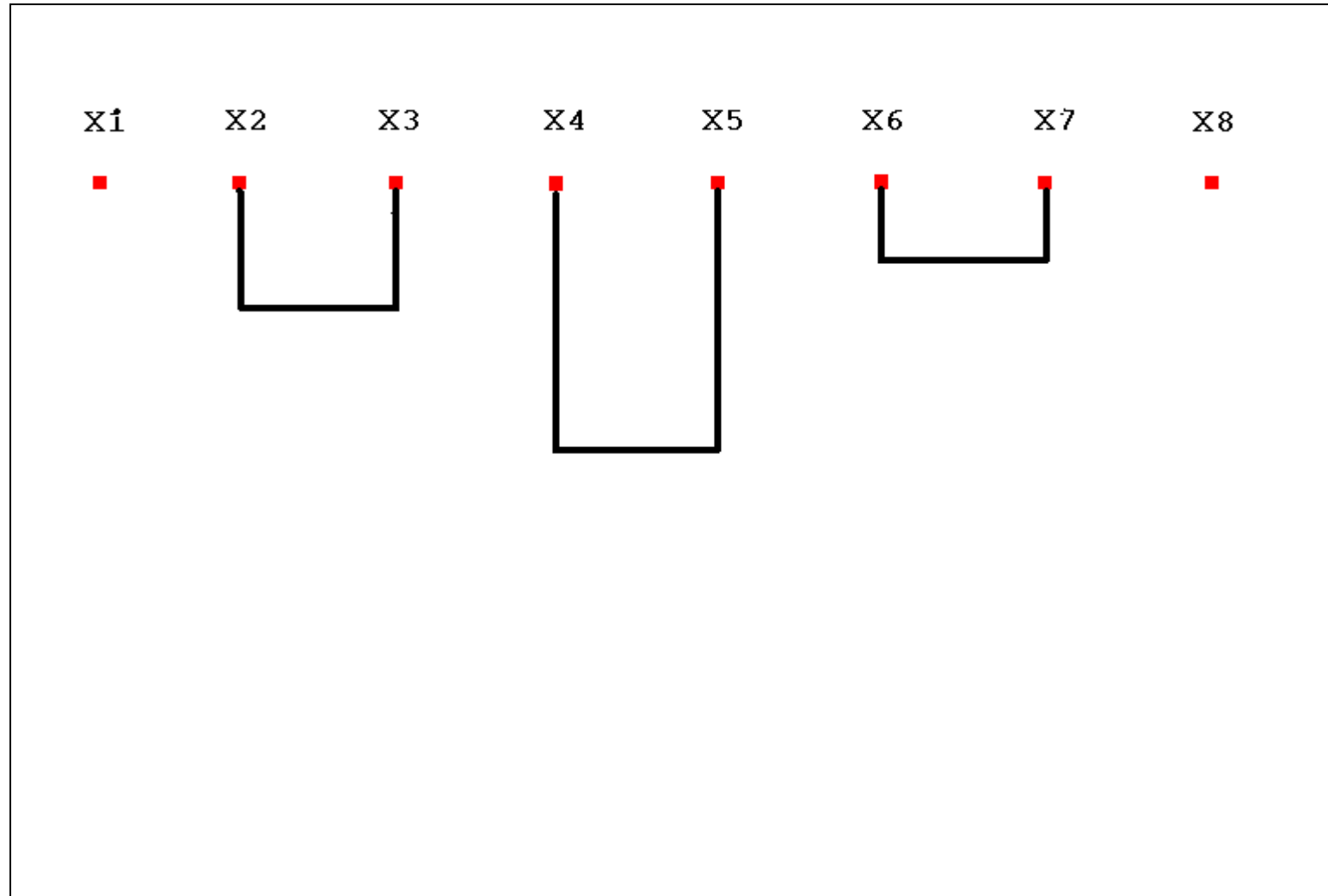


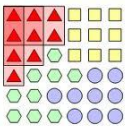
Nearest Neighbor, Level 3, $k = 6$ clusters.



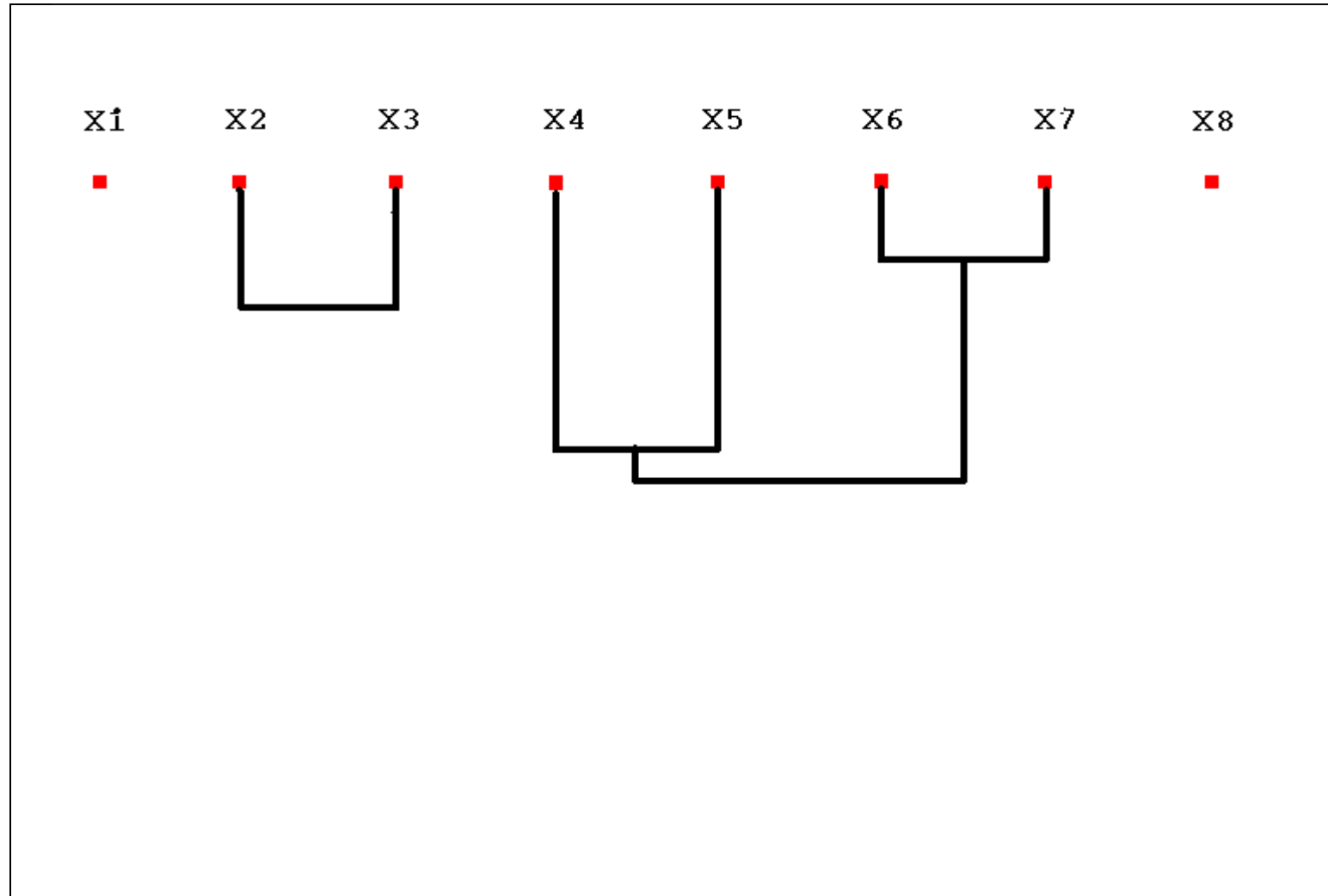


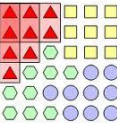
Nearest Neighbor, Level 4, $k = 5$ clusters.



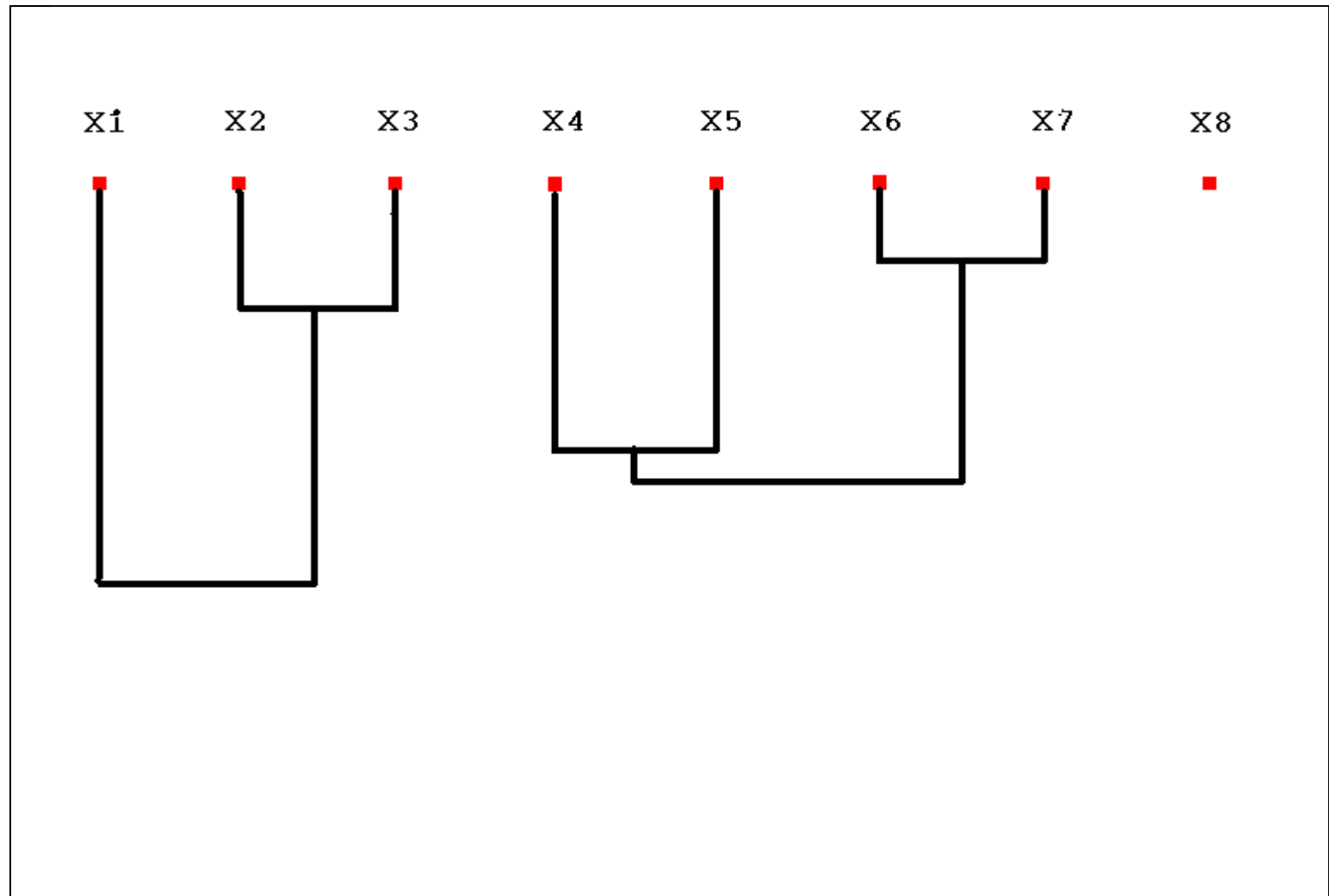


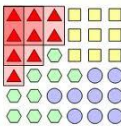
Nearest Neighbor, Level 5, $k = 4$ clusters.



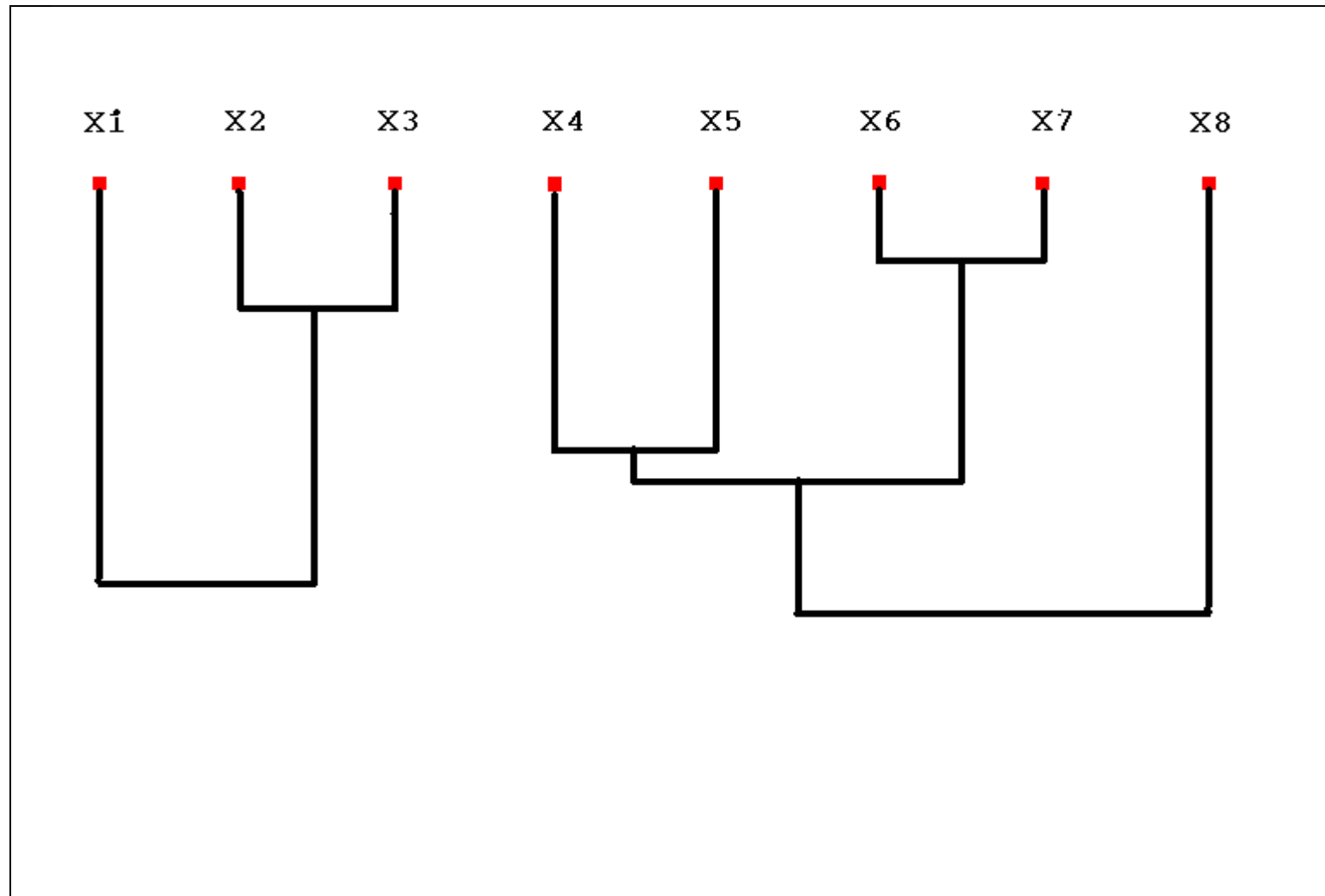


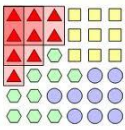
Nearest Neighbor, Level 6, k = 3 clusters.



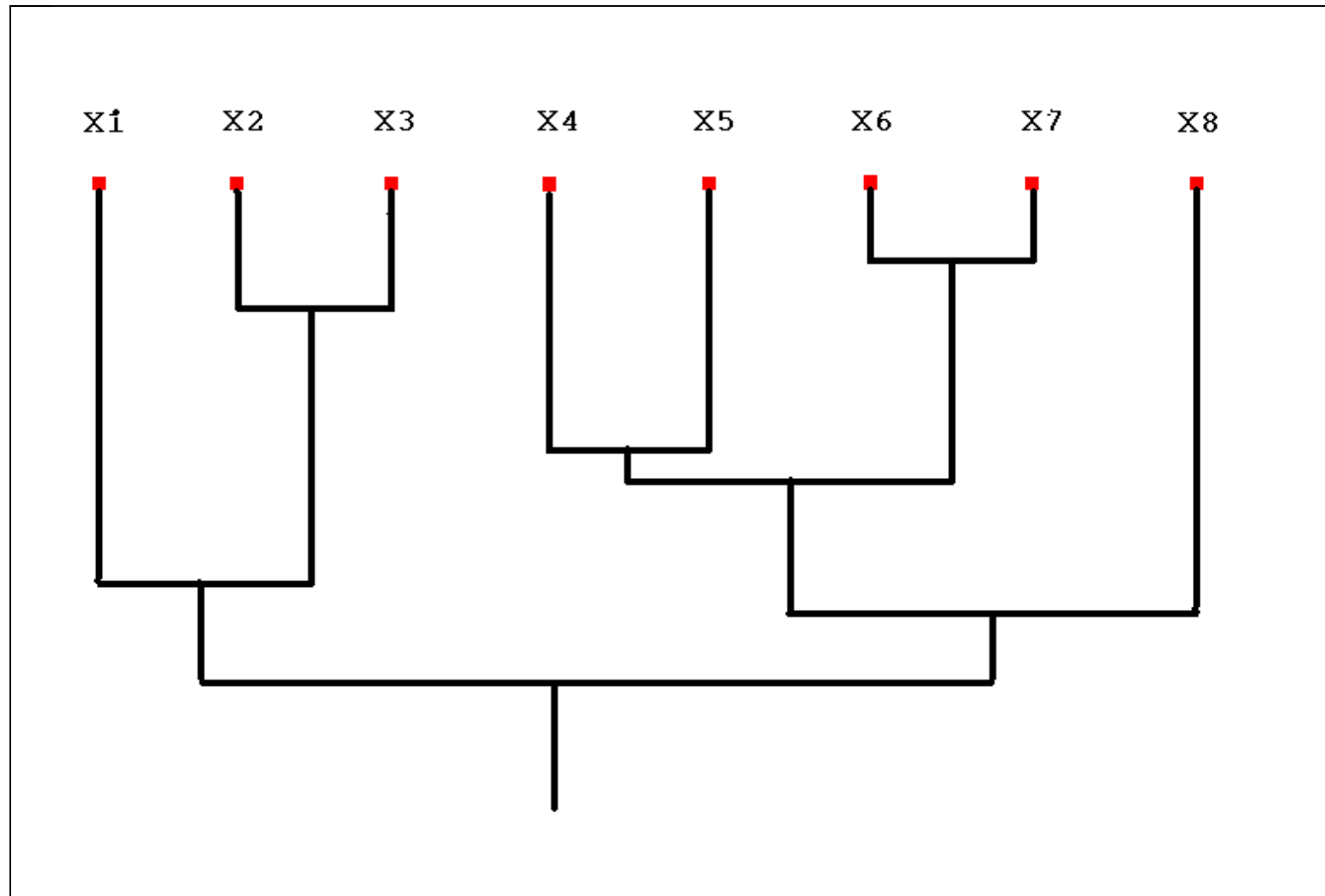


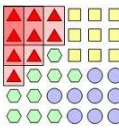
Nearest Neighbor, Level 7, $k = 2$ clusters.





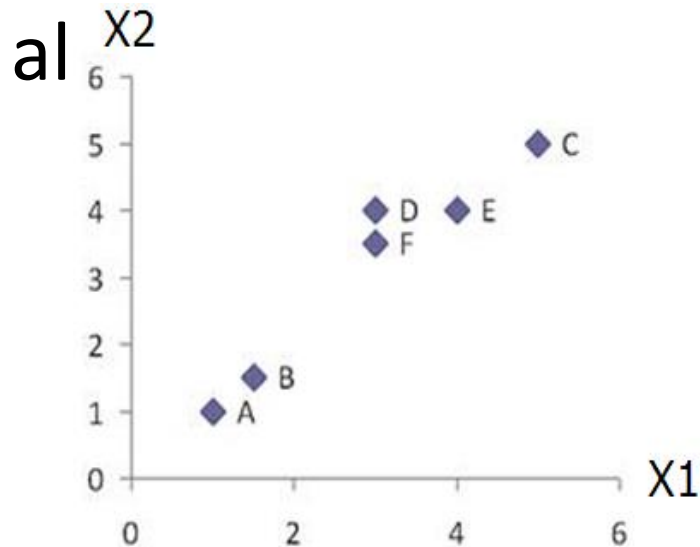
Nearest Neighbor, Level 8, $k = 1$ cluster.





Example and Demo

- Problem: clustering analysis with agglomerative



	X1	X2
A	1	1
B	1.5	1.5
C	5	5
D	3	4
E	4	4
F	3	3.5

data matrix

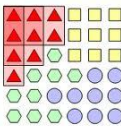
$$d_{AB} = \left((1-1.5)^2 + (1-1.5)^2 \right)^{\frac{1}{2}} = \sqrt{\frac{1}{2}} = 0.7071$$

$$d_{DF} = \left((3-3)^2 + (4-3.5)^2 \right)^{\frac{1}{2}} = 0.5$$

Euclidean distance

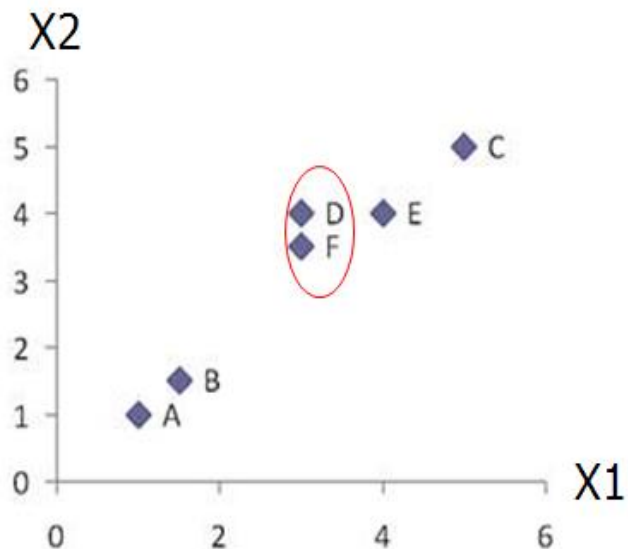
Dist	A	B	C	D	E	F
A	0.00	0.71	5.66	3.61	4.24	3.20
B	0.71	0.00	4.95	2.92	3.54	2.50
C	5.66	4.95	0.00	2.24	1.41	2.50
D	3.61	2.92	2.24	0.00	1.00	0.50
E	4.24	3.54	1.41	1.00	0.00	1.12
F	3.20	2.50	2.50	0.50	1.12	0.00

distance matrix



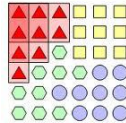
Example and Demo

- Merge two closest clusters



Dist	A	B	C	D	E	F
A	0.00	0.71	5.66	3.61	4.24	3.20
B	0.71	0.00	4.95	2.92	3.54	2.50
C	5.66	4.95	0.00	2.24	1.41	2.50
D	3.61	2.92	2.24	0.00	1.00	0.50
E	4.24	3.54	1.41	1.00	0.00	1.12
F	3.20	2.50	2.50	0.50	1.12	0.00

Dist	A	B	C	D, F	E
A	0.00	0.71	5.66	?	4.24
B	0.71	0.00	4.95	?	3.54
C	5.66	4.95	0.00	?	1.41
D, F	?	?	?	0.00	?
E	4.24	3.54	1.41	?	0.00



Example and Demo

- Update distance matrix

Dist	A	B	C	D	E	F
A	0.00	0.71	5.66	3.61	4.24	3.20
B	0.71	0.00	4.95	2.92	3.54	2.50
C	5.66	4.95	0.00	2.24	1.41	2.50
D	3.61	2.92	2.24	0.00	1.00	0.50
E	4.24	3.54	1.41	1.00	0.00	1.12
F	3.20	2.50	2.50	0.50	1.12	0.00

$$d_{(D,F) \rightarrow A} = \min(d_{DA}, d_{FA}) = \min(3.61, 3.20) = 3.20$$

$$d_{(D,F) \rightarrow B} = \min(d_{DB}, d_{FB}) = \min(2.92, 2.50) = 2.50$$

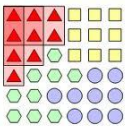
$$d_{(D,F) \rightarrow C} = \min(d_{DC}, d_{FC}) = \min(2.24, 2.50) = 2.24$$

$$d_{E \rightarrow (D,F)} = \min(d_{ED}, d_{EF}) = \min(1.00, 1.12) = 1.00$$

Dist	A	B	C	D, F	E
A	0.00	0.71	5.66	?	4.24
B	0.71	0.00	4.95	?	3.54
C	5.66	4.95	0.00	?	1.41
D, F	?	?	?	0.00	?
E	4.24	3.54	1.41	?	0.00

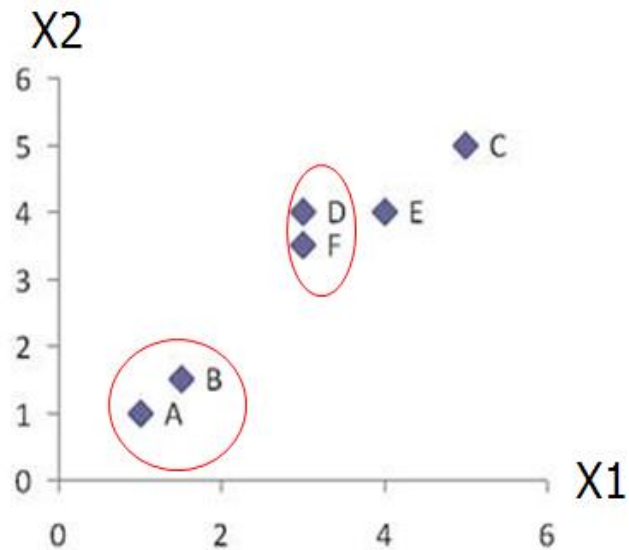
Min Distance (Single Linkage)

Dist	A	B	C	D, F	E
A	0.00	0.71	5.66	3.20	4.24
B	0.71	0.00	4.95	2.50	3.54
C	5.66	4.95	0.00	2.24	1.41
D, F	3.20	2.50	2.24	0.00	1.00
E	4.24	3.54	1.41	1.00	0.00



Example and Demo

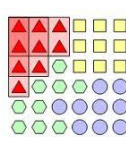
- Merge two closest clusters



Min Distance (Single Linkage)

Dist	A	B	C	D, F	E
A	0.00	0.71	5.66	3.20	4.24
B	0.71	0.00	4.95	2.50	3.54
C	5.66	4.95	0.00	2.24	1.41
D, F	3.20	2.50	2.24	0.00	1.00
E	4.24	3.54	1.41	1.00	0.00

Dist	A, B	C	(D, F)	E
A, B	0	?	?	?
C	?	0	2.24	1.41
(D, F)	?	2.24	0	1.00
E	?	1.41	1.00	0



Example and Demo

- Update distance matrix

Min Distance (Single Linkage)

Dist	A	B	C	D, F	E
A	0.00	0.71	5.66	3.20	4.24
B	0.71	0.00	4.95	2.50	3.54
C	5.66	4.95	0.00	2.24	1.41
D, F	3.20	2.50	2.24	0.00	1.00
E	4.24	3.54	1.41	1.00	0.00

$$d_{C \rightarrow (A,B)} = \min(d_{CA}, d_{CB}) = \min(5.66, 4.95) = 4.95$$

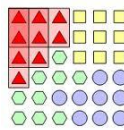
$$d_{(D,F) \rightarrow (A,B)} = \min(d_{DA}, d_{DB}, d_{FA}, d_{FB}) = \min(3.61, 2.92, 3.20, 2.50) = 2.50$$

$$d_{E \rightarrow (A,B)} = \min(d_{EA}, d_{EB}) = \min(4.24, 3.54) = 3.54$$

Dist	A,B	C	(D, F)	E
A,B	0	?	?	?
C	?	0	2.24	1.41
(D, F)	?	2.24	0	1.00
E	?	1.41	1.00	0

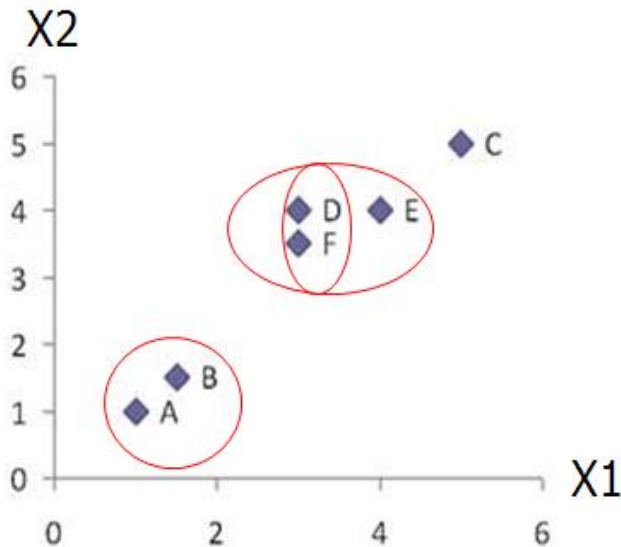
Min Distance (Single Linkage)

Dist	A,B	C	(D, F)	E
A,B	0	4.95	2.50	3.54
C	4.95	0	2.24	1.41
(D, F)	2.50	2.24	0	1.00
E	3.54	1.41	1.00	0



Example and Demo

- Merge two closest clusters/update distance matrix

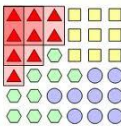


Min Distance (Single Linkage)

Dist	A,B	C	(D, F)	E
A,B	0	4.95	2.50	3.54
C	4.95	0	2.24	1.41
(D, F)	2.50	2.24	0	1.00
E	3.54	1.41	1.00	0

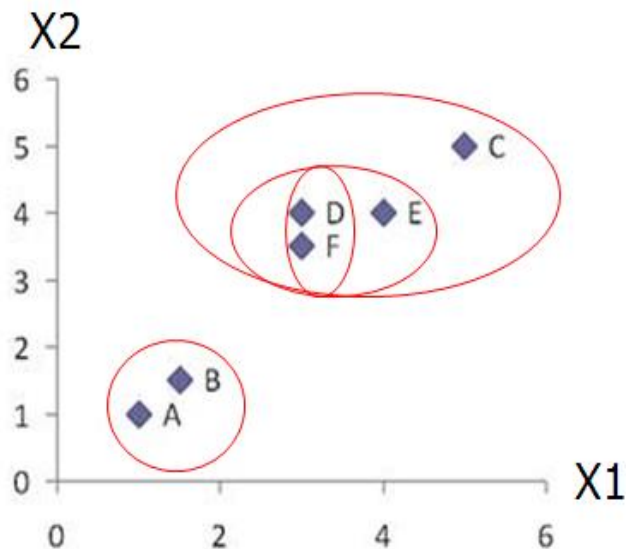
Min Distance (Single Linkage)

Dist	(A,B)	C	(D, F), E
(A,B)	0.00	4.95	2.50
C	4.95	0.00	1.41
(D, F), E	2.50	1.41	0.00



Example and Demo

- Merge two closest clusters/update distance matrix

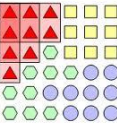


Min Distance (Single Linkage)

Dist	(A,B)	C	(D, F), E
(A,B)	0.00	4.95	2.50
C	4.95	0.00	1.41
(D, F), E	2.50	1.41	0.00

Min Distance (Single Linkage)

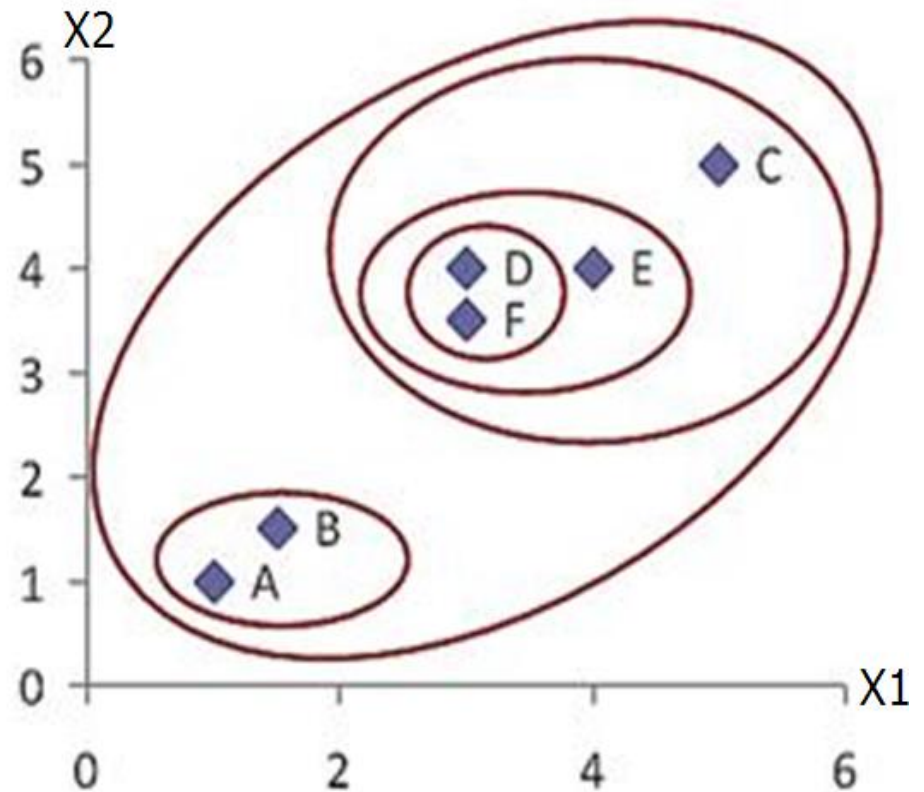
Dist	(A,B)	((D, F), E), C
(A,B)	0.00	2.50
((D, F), E), C	2.50	0.00

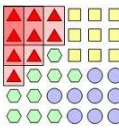


Example and Demo

- Final result (meeting termination condition)

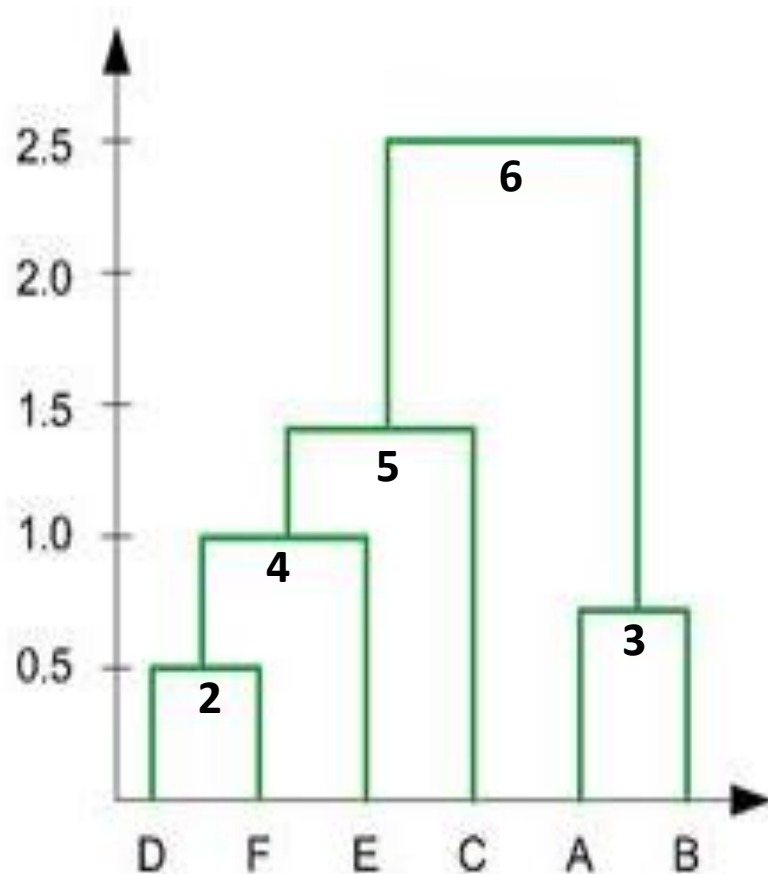
	X1	X2
A	1	1
B	1.5	1.5
C	5	5
D	3	4
E	4	4
F	3	3.5





Example and Demo

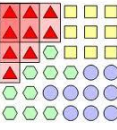
- **Dendrogram tree** representation



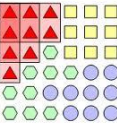
1. In the beginning we have 6 clusters: A, B, C, D, E and F
2. We merge cluster D and F into cluster (D, F) at distance 0.50
3. We merge cluster A and cluster B into (A, B) at distance 0.71
4. We merge cluster E and (D, F) into ((D, F), E) at distance 1.00
5. We merge cluster ((D, F), E) and C into (((D, F), E), C) at distance 1.41
6. We merge cluster (((D, F), E), C) and (A, B) into ((((D, F), E), C), (A, B)) at distance 2.50
7. The last cluster contain all the objects, thus conclude the computation



Cluster Analysis

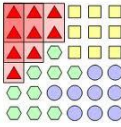


- What is Cluster Analysis?
- Types of Data in Cluster Analysis
- A Categorization of Major Clustering Methods
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Grid-Based Methods
- Model-Based Clustering Methods
- Summary



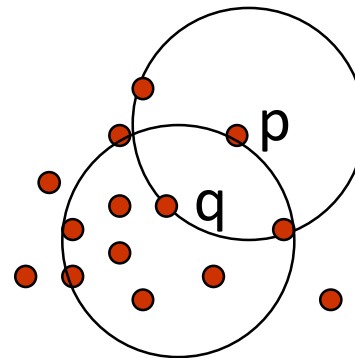
Density-Based Clustering Methods

- 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)

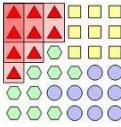


Density-Based Clustering: Background

- Neighborhood of point p = all points within distance Eps from p :
 - $N_{Eps}(p) = \{q \mid \text{dist}(p, q) \leq Eps\}$
- Two parameters:
 - **Eps**: Maximum radius of the neighbourhood
 - **MinPts**: Minimum number of points in an Eps-neighbourhood of that point
- If the number of points in the Eps-neighborhood of p is at least **MinPts**, then p is called a **core object**.
- **Directly density-reachable**: A point p is directly density-reachable from a point q wrt. **Eps**, **MinPts** if
 - 1) p belongs to $N_{Eps}(q)$
 - 2) core point condition:
 $|N_{Eps}(q)| \geq \text{MinPts}$



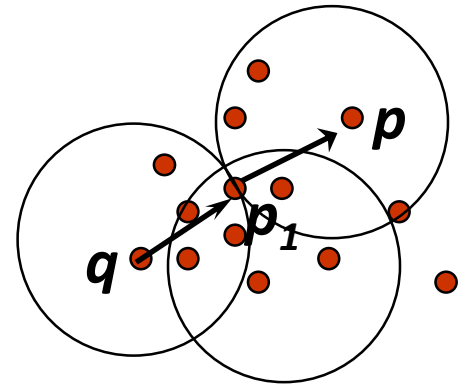
MinPts = 5
Eps = 1 cm



Density-Based Clustering: Background (II)

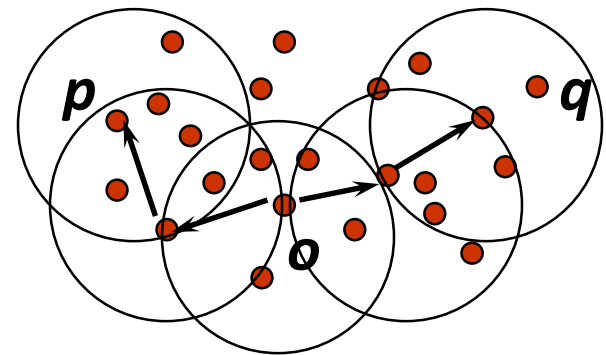
- **Density-reachable:**

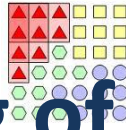
- A point p is density-reachable from a point q wrt. Eps , $MinPts$ if there is a chain of points $p_1, \dots, p_n, p_1 = q, p_n = p$ such that p_{i+1} is directly density-reachable from p_i



- **Density-connected**

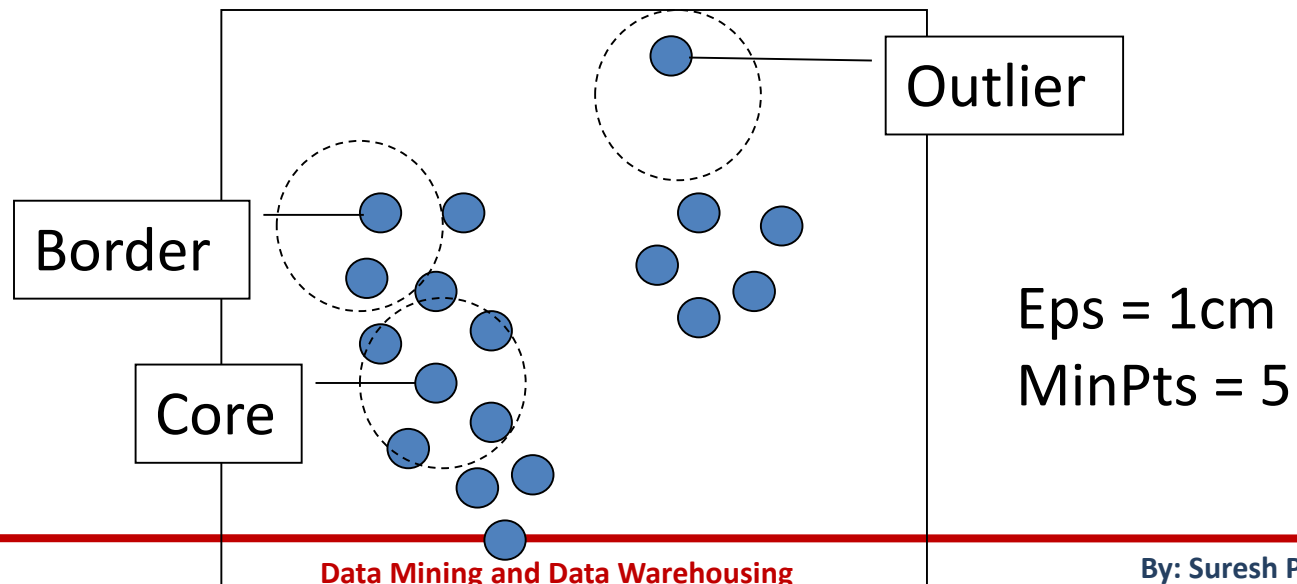
- A point p is density-connected to a point q wrt. Eps , $MinPts$ if there is a point o such that both, p and q are density-reachable from o wrt. Eps and $MinPts$.

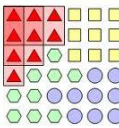




DBSCAN: Density Based Spatial Clustering of Applications with Noise

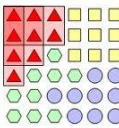
- Relies on a *density-based* notion of cluster: A *cluster* is defined as a maximal set of density-connected points
- Discovers clusters of arbitrary shape in spatial databases with noise





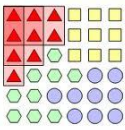
DBSCAN: The Algorithm

- Arbitrary select a point p
- Retrieve all points density-reachable from p wrt ***Eps*** and ***MinPts***.
- If p is a core point, a cluster is formed.
- If p is a border point, no points are density-reachable from p and DBSCAN visits the next point of the database.
- Continue the process until all of the points have been processed.



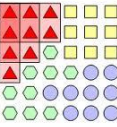
Cluster Analysis

- What is Cluster Analysis?
- Types of Data in Cluster Analysis
- A Categorization of Major Clustering Methods
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- **Grid-Based Methods**
- Model-Based Clustering Methods
- Summary



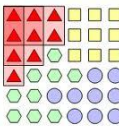
Grid-Based Clustering Method

- Using multi-resolution grid data structure
- Several interesting methods
 - **STING** (a STatistical INformation Grid approach) by Wang, Yang and Muntz (1997)
 - **WaveCluster** by Sheikholeslami, Chatterjee, and Zhang (VLDB'98)
 - A multi-resolution clustering approach using wavelet method
 - **CLIQUE**: Agrawal, et al. (SIGMOD'98)



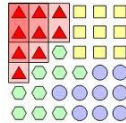
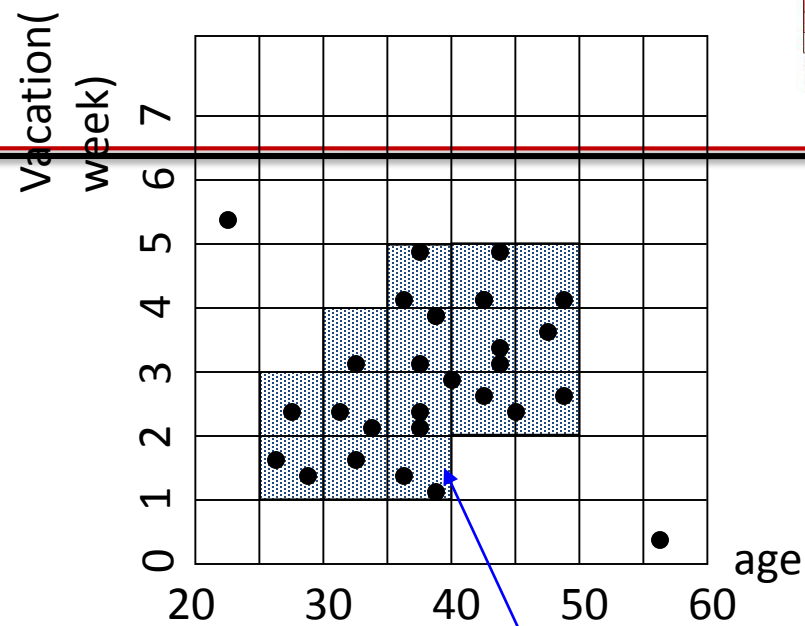
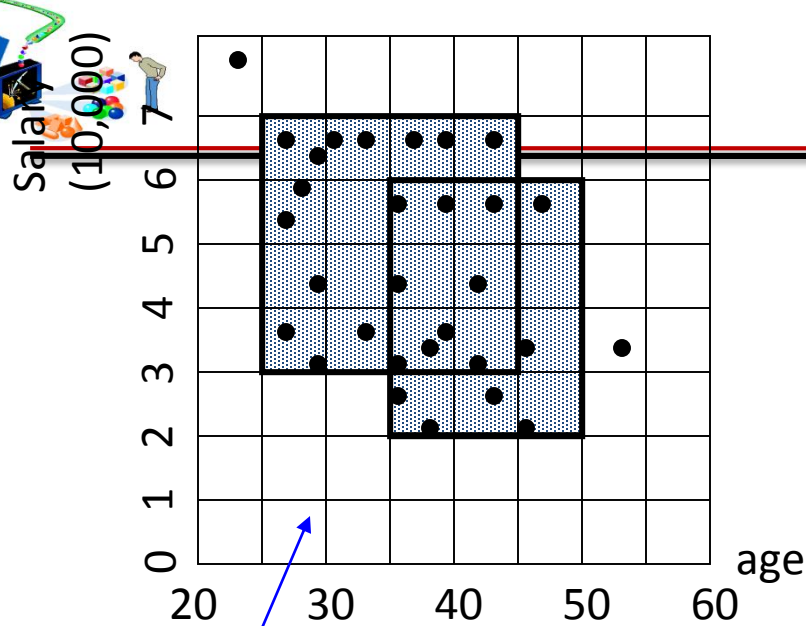
CLIQUE (Clustering In QUES)

- Agrawal, Gehrke, Gunopulos, Raghavan (SIGMOD'98).
- Automatically identifying subspaces of a high dimensional data space that allow better clustering than original space
- CLIQUE can be considered as both density-based and grid-based
 - It partitions each dimension into the same number of equal length interval
 - It partitions an m-dimensional data space into non-overlapping rectangular units
 - A unit is dense if the fraction of total data points contained in the unit exceeds the input model parameter
 - A cluster is a maximal set of connected dense units within a subspace



CLIQUE: The Major Steps

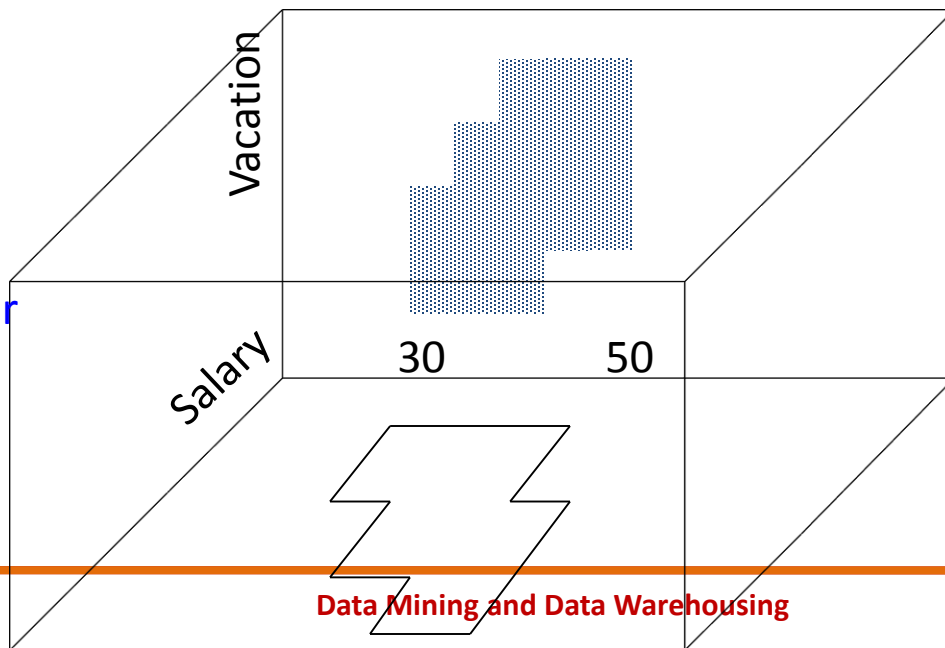
- Partition the data space and find the number of points that lie inside each cell of the partition.
- Identify the subspaces that contain clusters using the Apriori principle
- Identify clusters:
 - Determine dense units in all subspaces of interests
 - Determine connected dense units in all subspaces of interests.
- Generate minimal description for the clusters
 - Determine maximal regions that cover a cluster of connected dense units for each cluster
 - Determination of minimal cover for each cluster
- CLIQUE can find **projected clusters** in subspaces of the dimensional space



$\tau = 3$

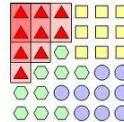
projected cluster
in (salary, age)
subspace

projected cluster
in (vacation,
age) subspace





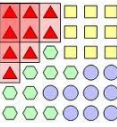
Cluster Analysis



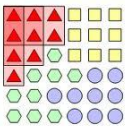
- What is Cluster Analysis?
- Types of Data in Cluster Analysis
- A Categorization of Major Clustering Methods
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Grid-Based Methods
- **Model-Based Clustering Methods**
- Summary



Model based clustering



- Assume data generated from K probability distributions
- Typically Gaussian distribution Soft or probabilistic version of K-means clustering
- Need to find distribution parameters.
- EM Algorithm



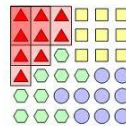
EM Algorithm

- Initialize K cluster centers
- Iterate between two steps
 - **E**xpectation step: assign points to clusters

$$P(d_i \in c_k) = \frac{w_k \Pr(d_i | c_k)}{\sum_j w_j \Pr(d_i | c_j)}$$
$$w_k = \frac{\sum_i \Pr(d_i \in c_k)}{N}$$

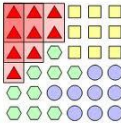
- **M**aximation step: estimate model parameters

$$\mu_k = \frac{1}{m} \sum_{i=1}^m \frac{d_i P(d_i \in c_k)}{\sum_j P(d_i \in c_j)}$$



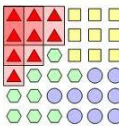
Cluster Analysis

- What is Cluster Analysis?
- Types of Data in Cluster Analysis
- A Categorization of Major Clustering Methods
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Grid-Based Methods
- Model-Based Clustering Methods
- Summary

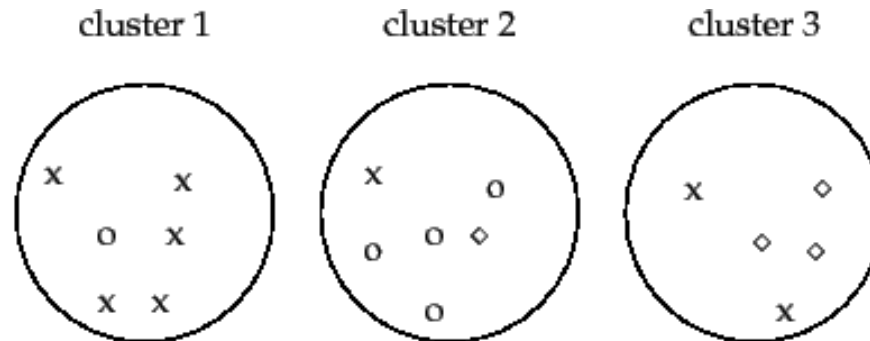


Validation Criteria : Clustering

- Purity
- *Rand index*
- *F measure*



Clustering: Purity

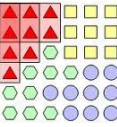


► **Figure 16.1** Purity as an external evaluation criterion for cluster quality. Majority class and number of members of the majority class for the three clusters are: x, 5 (cluster 1); o, 4 (cluster 2); and \diamond , 3 (cluster 3). Purity is $(1/17) \times (5 + 4 + 3) \approx 0.71$.

$$\text{purity}(\Omega, \mathbb{C}) = \frac{1}{N} \sum_k \max_j |\omega_k \cap c_j|$$

where $\Omega = \{\omega_1, \omega_2, \dots, \omega_K\}$ is the set of cluster and $\mathbb{C} = \{c_1, c_2, \dots, c_J\}$ is the set of classes

High purity is easy to achieve when the number of clusters is large - in particular, purity is 1 if each document gets its own cluster. Thus, we cannot use purity to trade off the quality of the clustering against the number of clusters.



Clustering: Rand Index

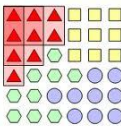
A true positive (TP) decision assigns two similar documents to the same cluster, a true negative (TN) decision assigns two dissimilar documents to different clusters. There are two types of errors we can commit. A (FP) decision assigns two dissimilar documents to the same cluster. A (FN) decision assigns two similar documents to different clusters. The *Rand index* () measures the percentage of decisions that are correct.

$$RI = \frac{TP + TN}{TP + FP + FN + TN}$$

$$TP + FP = \binom{6}{2} + \binom{6}{2} + \binom{5}{2} = 40$$

$$TP = \binom{5}{2} + \binom{4}{2} + \binom{3}{2} + \binom{2}{2} = 20$$

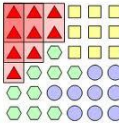
$$FP = 40 - 20 = 20$$



Clustering: Rand Index

	Same cluster	Different clusters
Same class	$TP = 20$ _____	$FN = 24$ _____
Different classes	$FP = 20$ _____	$TN = 72$ _____

RI is then $(20 + 72) / (20 + 20 + 24 + 72) \approx 0.68$.



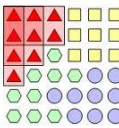
Clustering: F-Measure

	Same cluster	Different clusters
Same class	<u>TP = 20</u>	<u>FN = 24</u>
Different classes	<u>FP = 20</u>	<u>TN = 72</u>

$$P = \frac{TP}{TP + FP} \quad R = \frac{TP}{TP + FN} \quad F_{\beta} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

$$P = 20/40 = 0.5 \quad R = 20/44 \approx 0.455 \quad \underline{F_1 \approx 0.48} \quad \beta = 1$$
$$\underline{F_5 \approx 0.456} \quad \beta = 5$$

In information retrieval, evaluating clustering with \underline{F} has the advantage that the measure is already familiar to the research community.



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
- **Outlier detection** and analysis are very useful for fraud detection, etc. and can be performed by statistical, distance-based or deviation-based approaches

