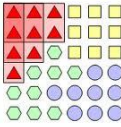


Data Mining and Data Warehousing

Unit 8

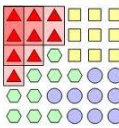
Classification and Prediction

Suresh Pokharel



What is classification?

- is a data mining technique used to predict the category of categorical data by building a model based on some predictor variables (to classify data).
- Predictor variable/attribute is called **class label attribute (predefined class)**



What is classification?

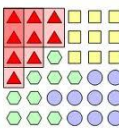
It is a **two-step** process

1. Model Construction (learning step or training phase)

- build a model to explain the target concept
- model is represented as classification rules, decision trees, or mathematical formulae.

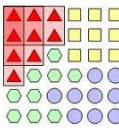
2. Model Usage

- is used for classifying future or unknown cases
- estimate the accuracy of the model

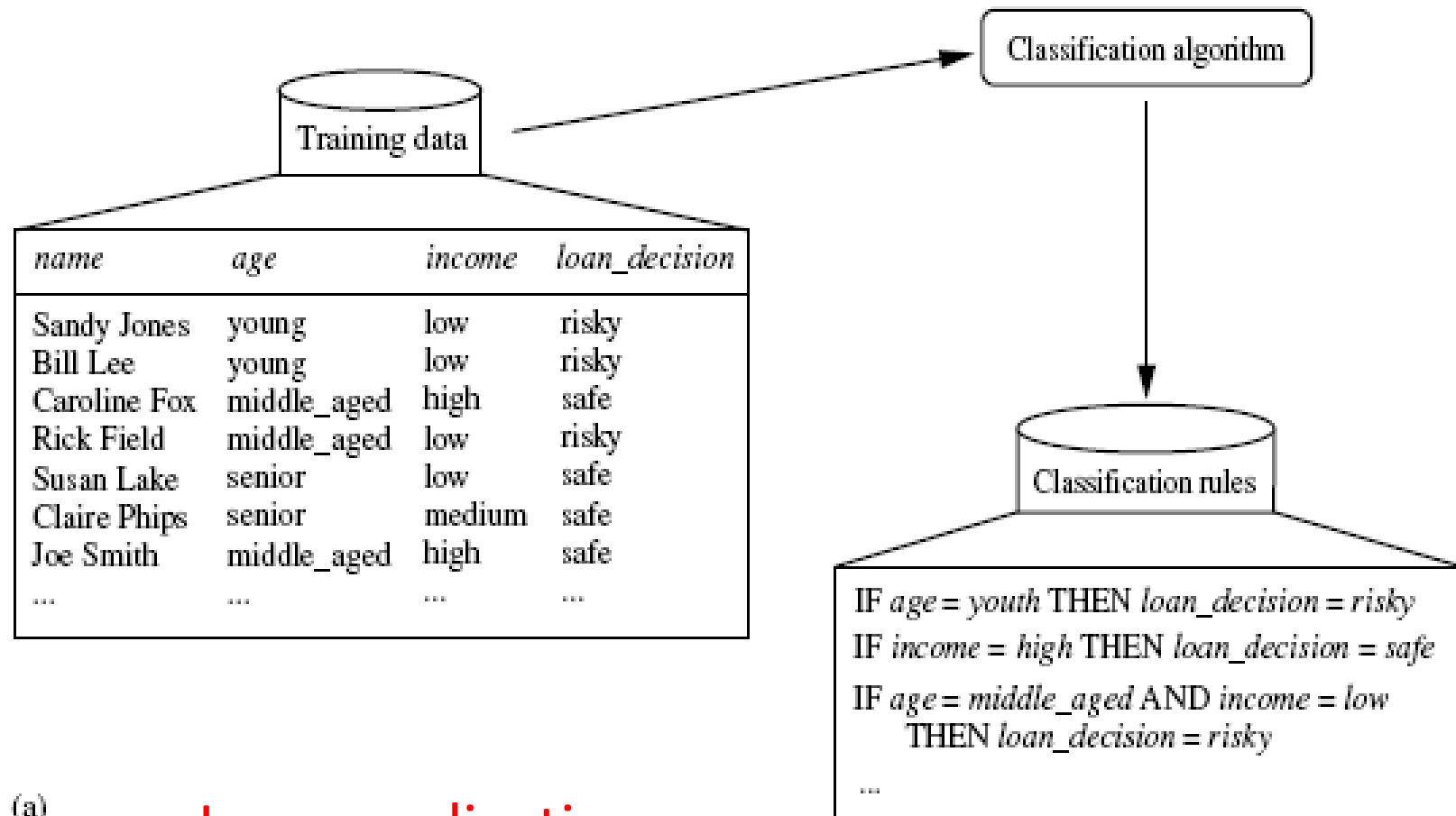


Example

<i>RID</i>	<i>age</i>	<i>income</i>	<i>student</i>	<i>credit_rating</i>	<i>Class: buys_computer</i>
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

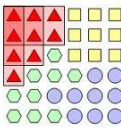


Step 1

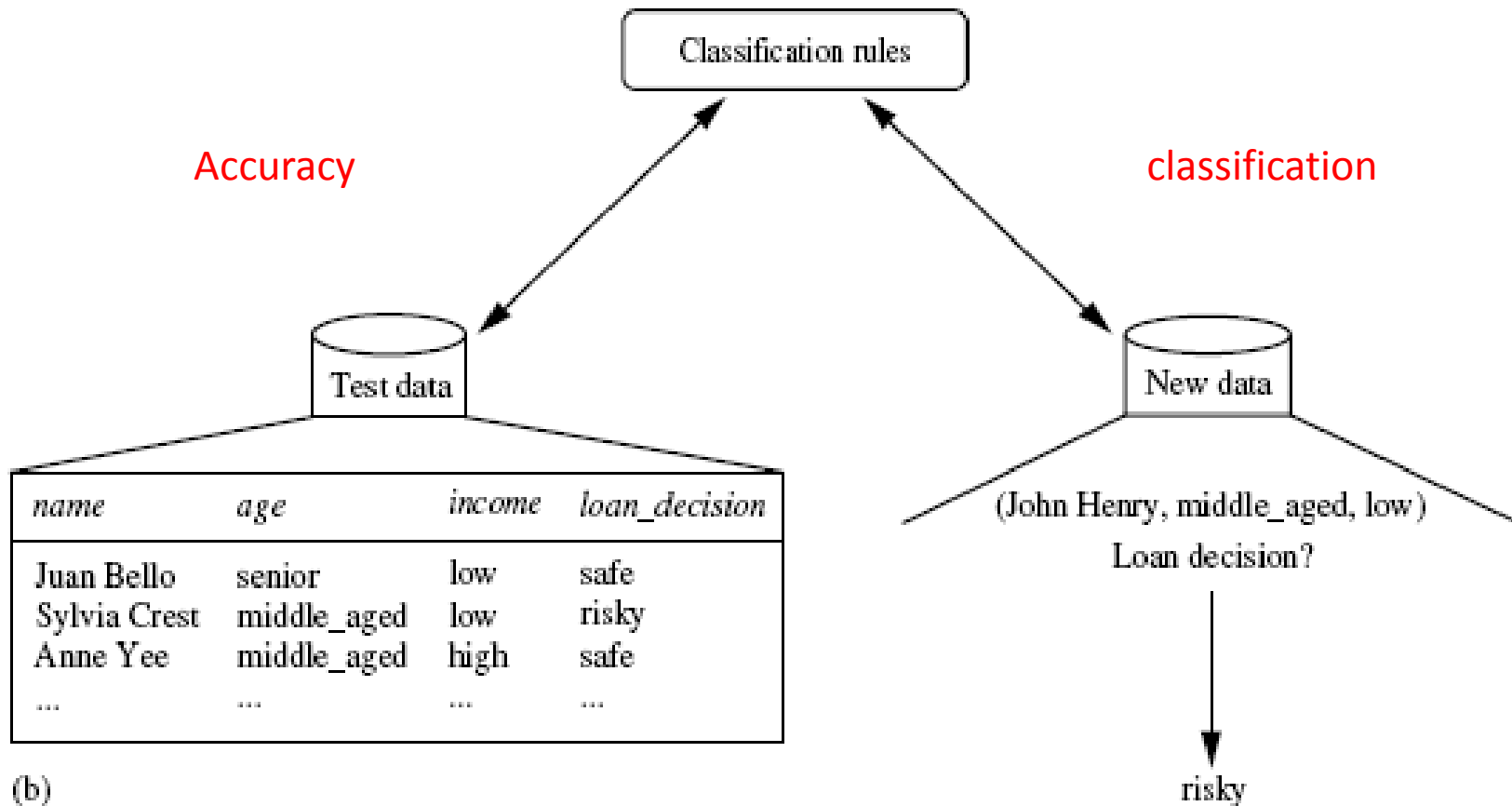


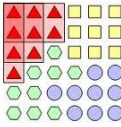
(a)

Loan application



Step 2



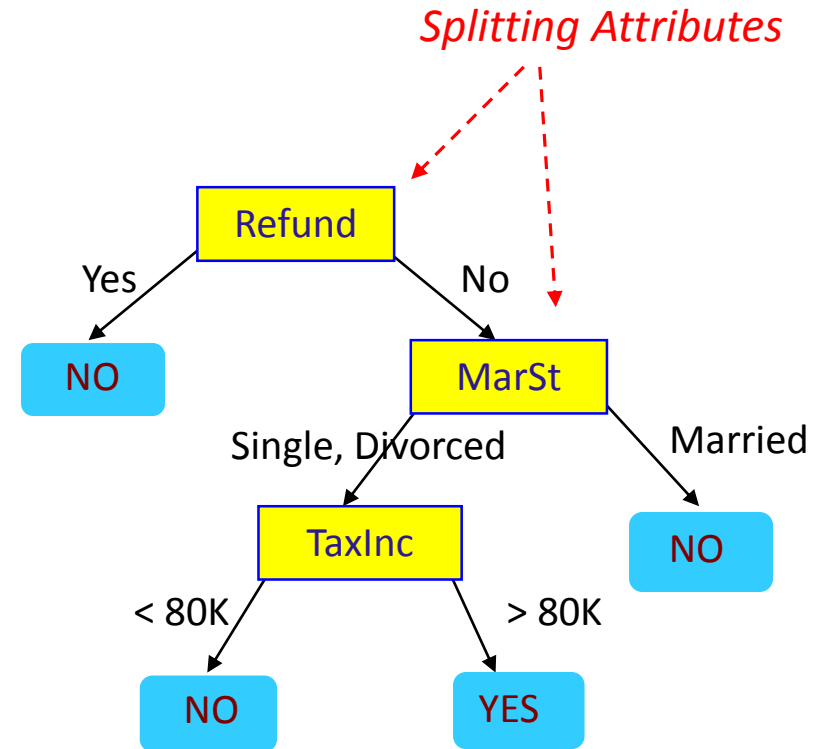


Example of a Decision Tree

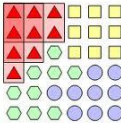
categorical
categorical
continuous
class

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Training Data



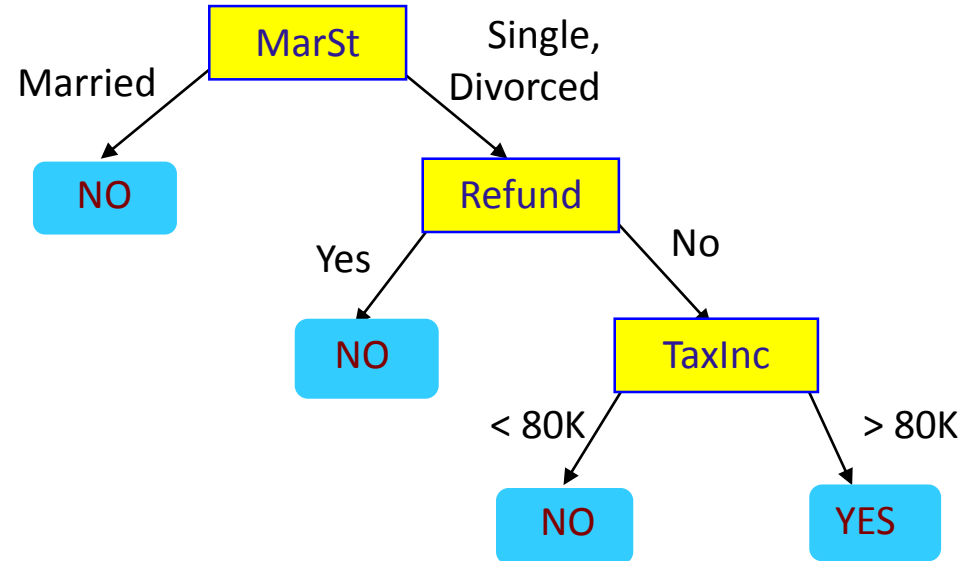
Model: Decision Tree



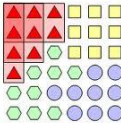
Another Example of Decision Tree

categorical
categorical
continuous
class

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



There could be more than one tree that fits the same data!



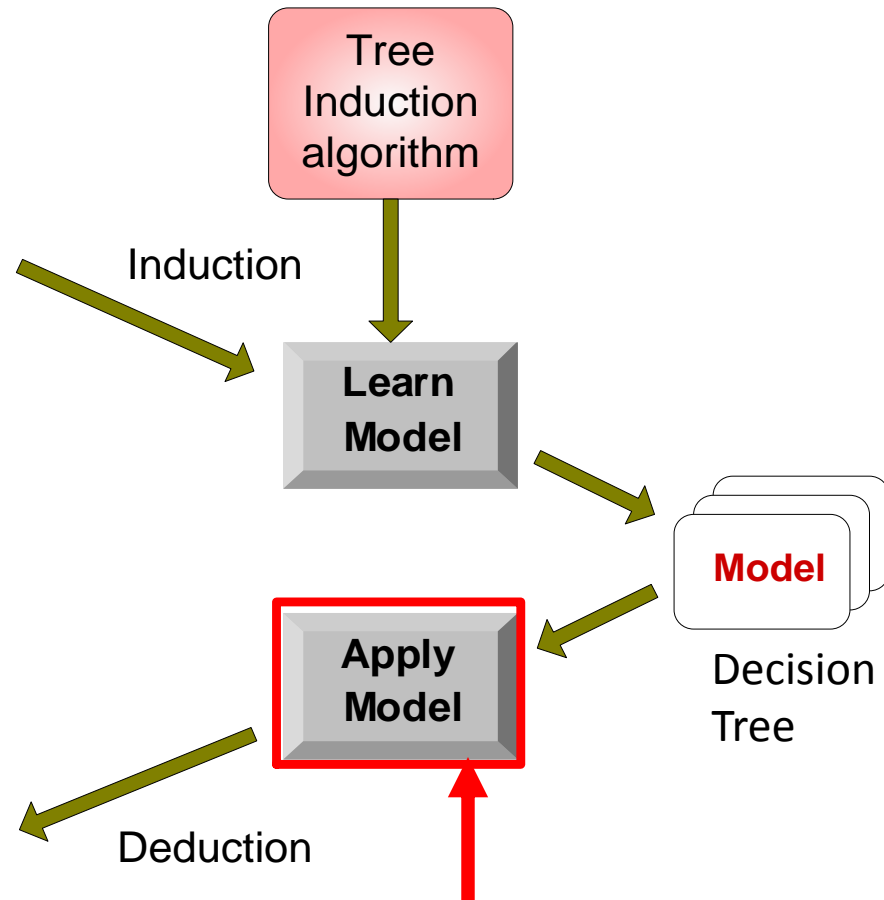
Decision Tree Classification Task

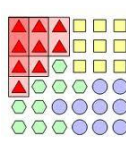
Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

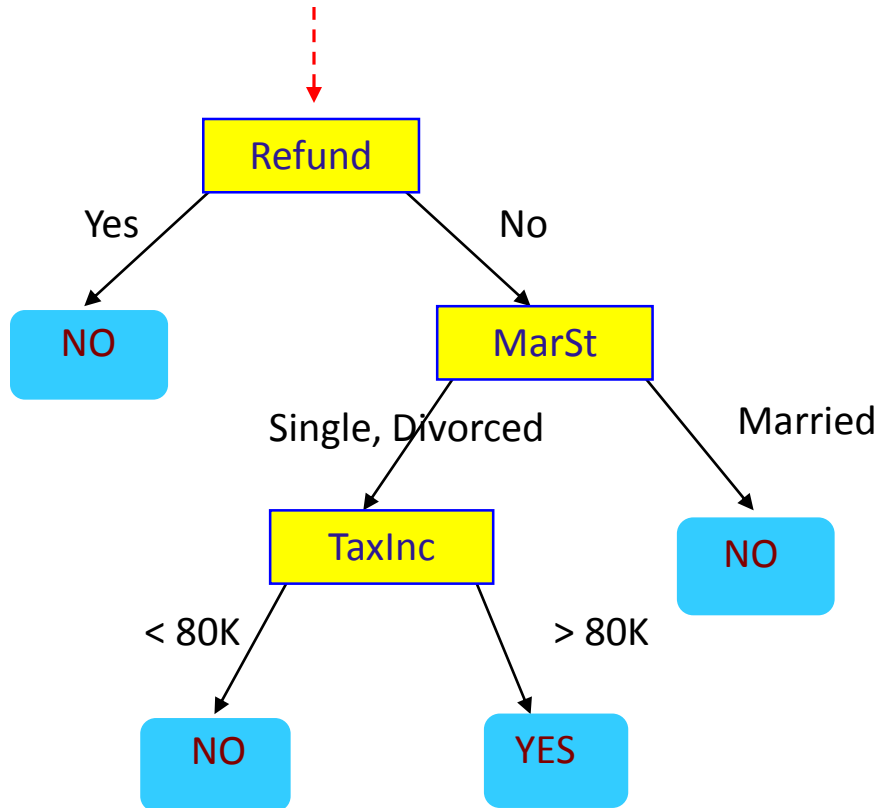
Test Set





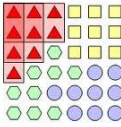
Apply Model to Test Data

Start from the root of tree.



Test Data

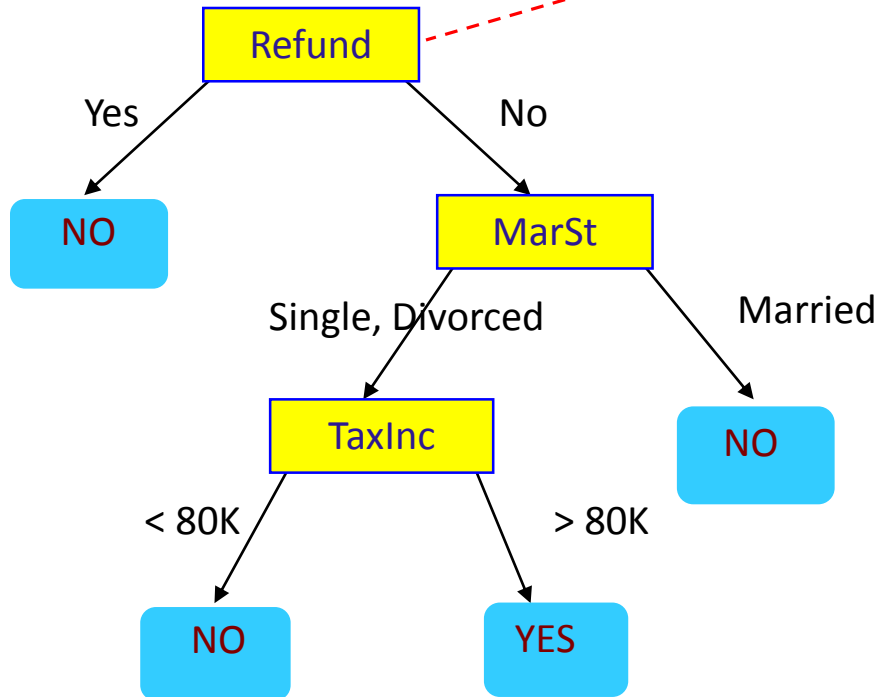
Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

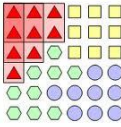


Apply Model to Test Data

Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

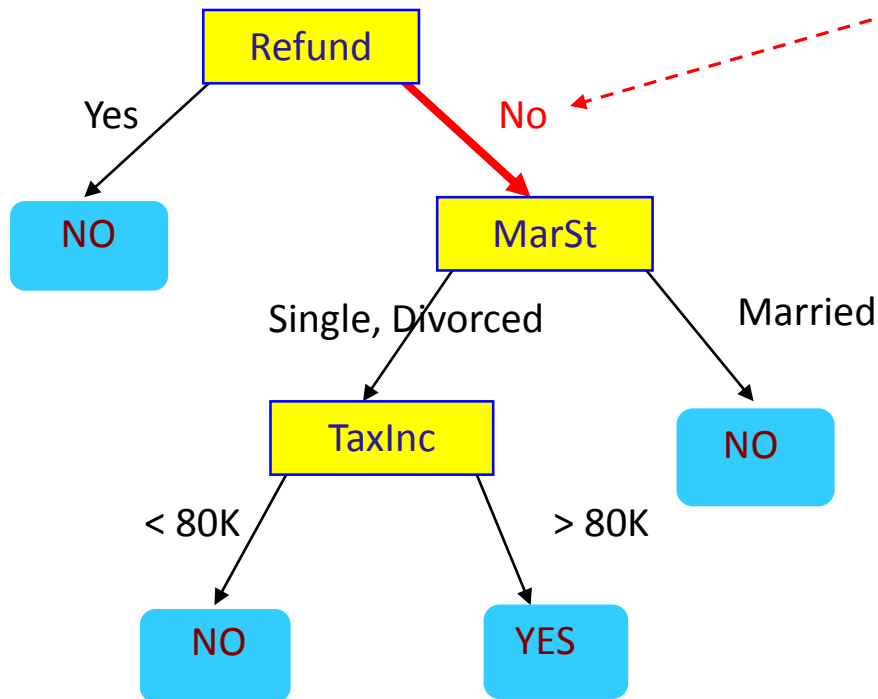


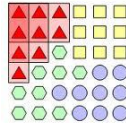


Apply Model to Test Data

Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

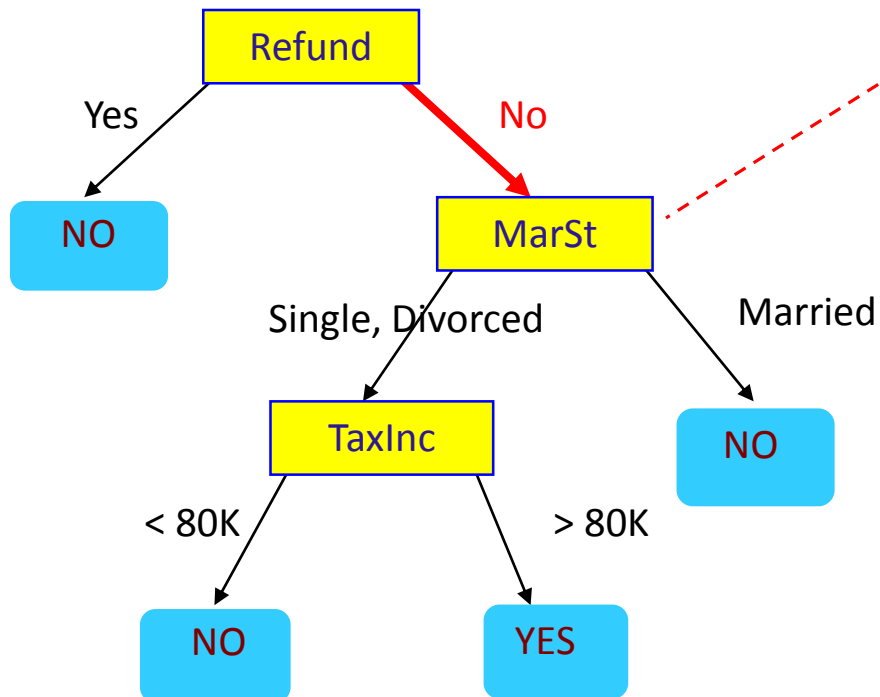


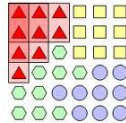


Apply Model to Test Data

Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

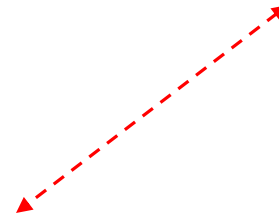
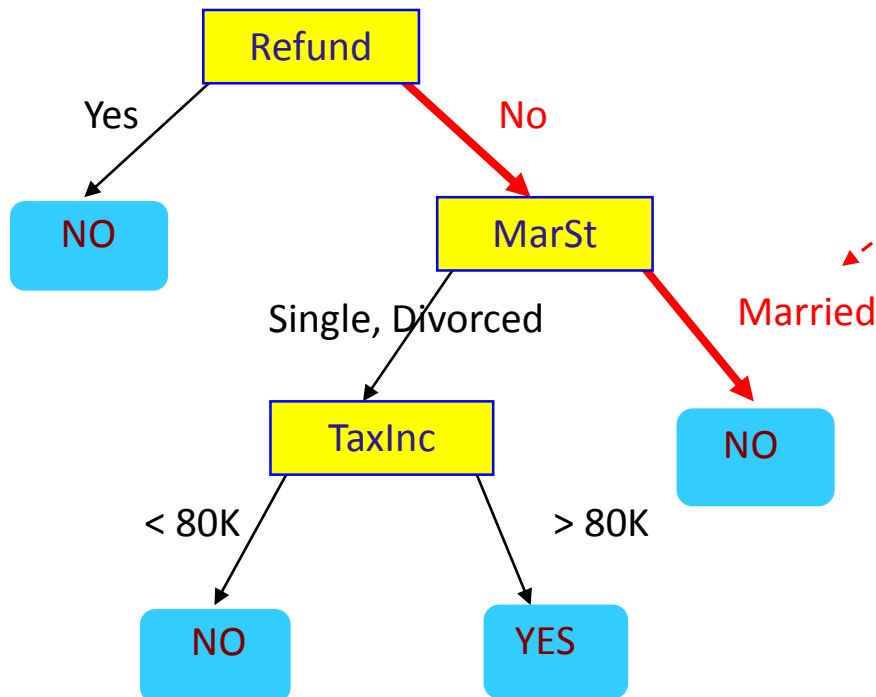


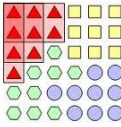


Apply Model to Test Data

Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

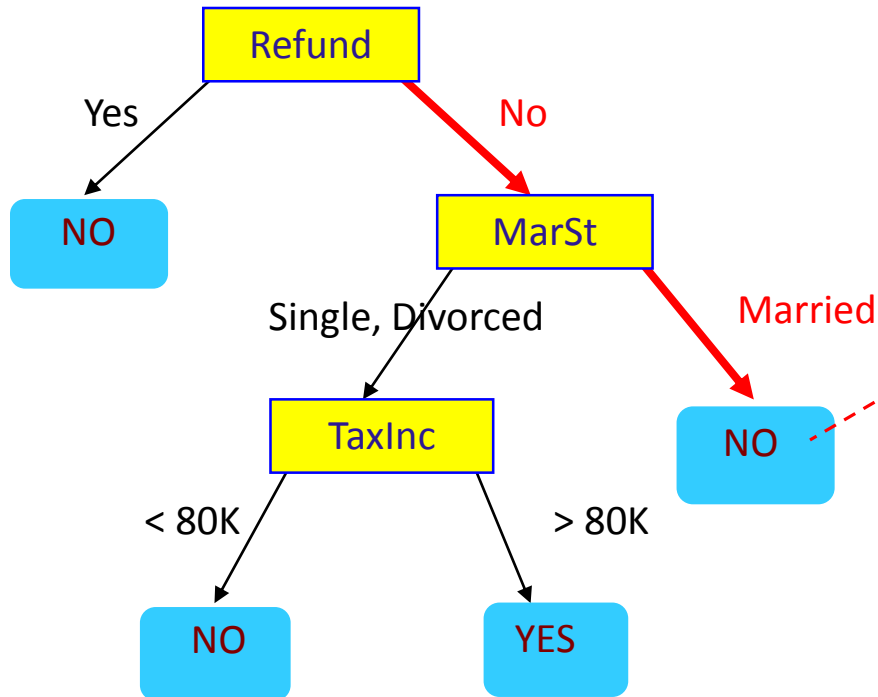




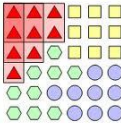
Apply Model to Test Data

Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



Assign Cheat to "No"



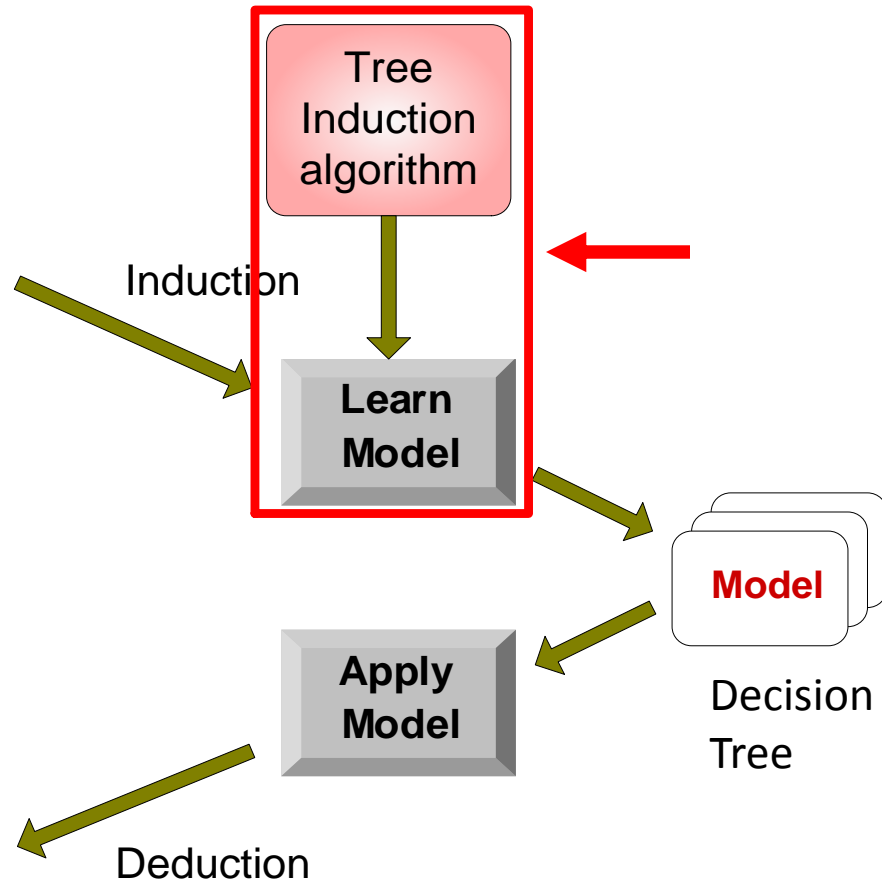
Decision Tree Classification Task

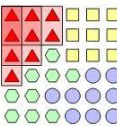
Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



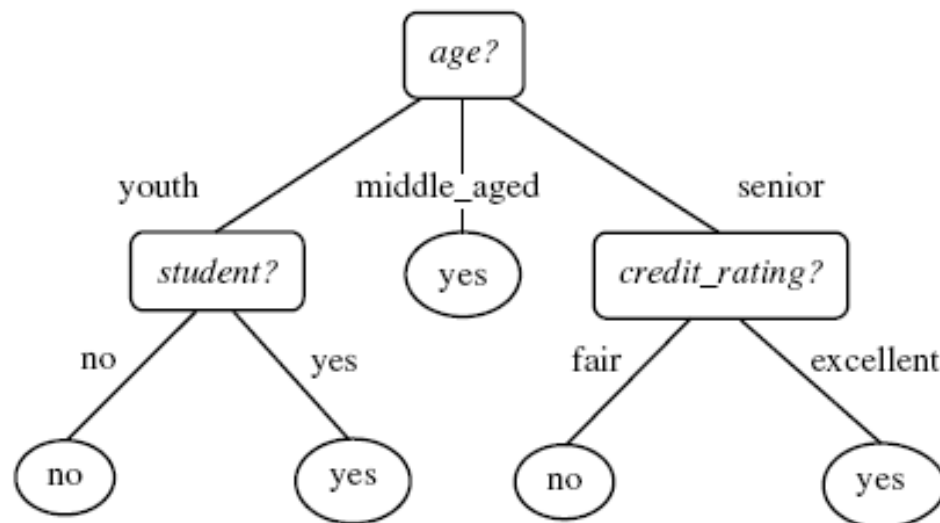


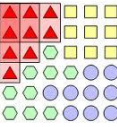
Classification by decision tree induction

What is decision tree?

- flow chart like tree structure
- internal node denotes a test on an attribute
- each branch represents an outcome of the test
- each leaf node holds a class label

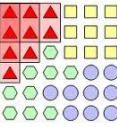
Decision tree for
the concept **buys
computer**





Why Decision Tree?

- Construction does not require any domain knowledge
- Can handle high dimensional data
- Learning step is simple and fast
- In general have good accuracy
- People are able to understand decision tree models after a brief explanation



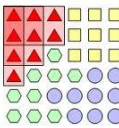
Issue regarding classification & Prediction

- Data Preparation

data cleaning, relevant analysis, data transformation and reduction

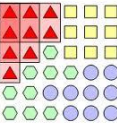
- Comparing classification & prediction method

Accuracy, speed, Robustness, scalability, interpretability



Supervised & Unsupervised Learning

- **Supervised Learning (Classification)**
 - Supervision: The training data are accompanied by labels indicating the class of the observations
 - New data is classified based on the training set
- **Unsupervised Learning (Clustering)**
 - The class labels of training data is unknown
 - Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

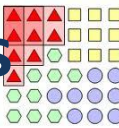


Some Decision Tree Systems

- ID3 (Quinlan 79)
- CART (Brieman et al. 84)
- Assistant (Cestnik et al. 87)
- C4.5 (Quinlan 93)
- See5 (Quinlan 97)
- ...
- Orange (Demšar, Zupan 98-03)

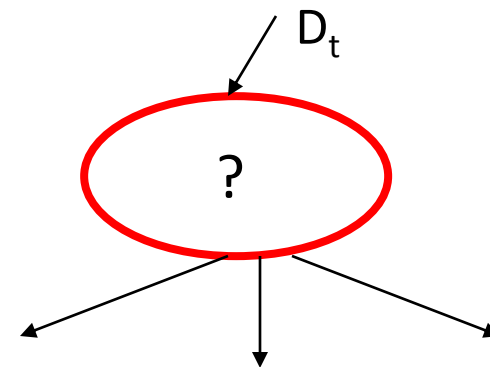


Decision Tree Induction- General Structure of Hunt's Algorithm

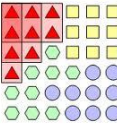


- Let D_t be the set of training records that reach a node t
- General Procedure:
 - If D_t contains records that belong the same class y_t , then t is a leaf node labeled as y_t
 - If D_t is an empty set, then t is a leaf node labeled by the default class, y_d
 - If D_t contains records that belong to more than one class, use an attribute test to split the data into smaller subsets. Recursively apply the procedure to each subset.

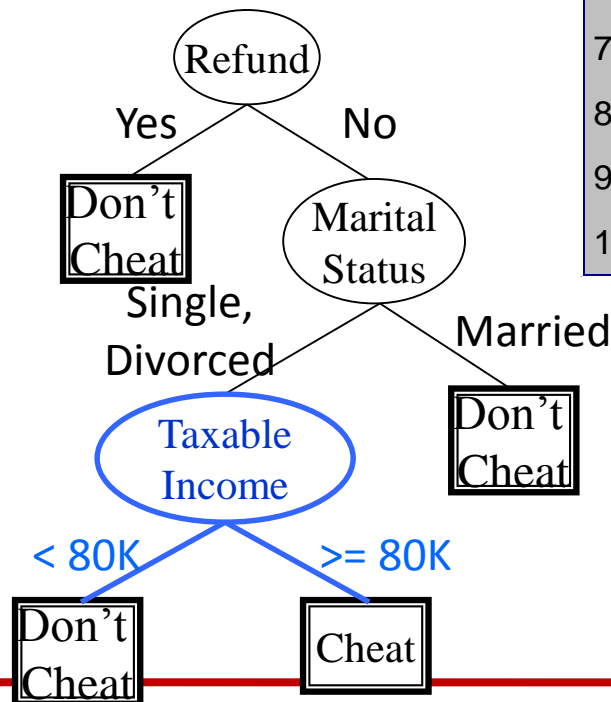
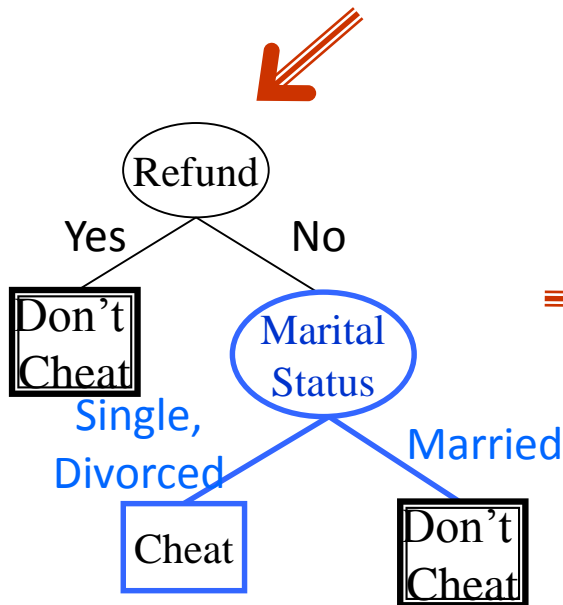
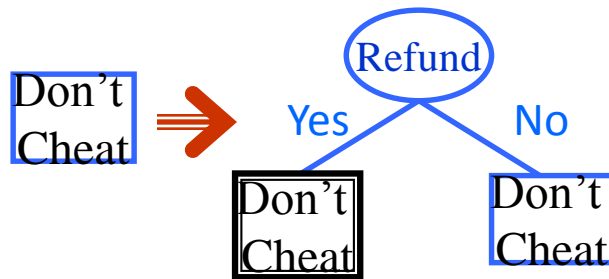
<i>Tid</i>	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

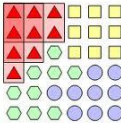


Decision Tree Induction- Hunt's Algorithm



Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes





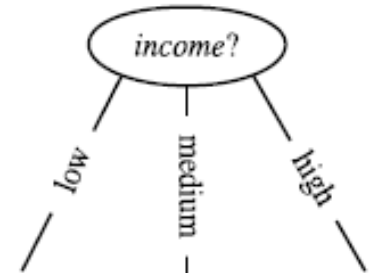
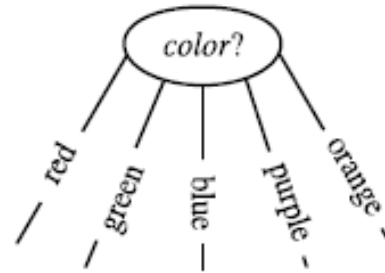
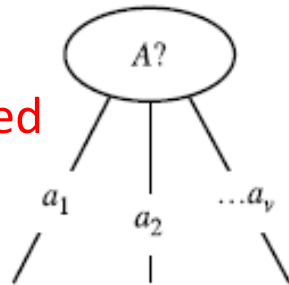
Example

Partitioning Scenarios

Examples

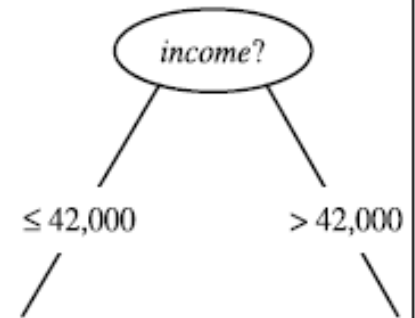
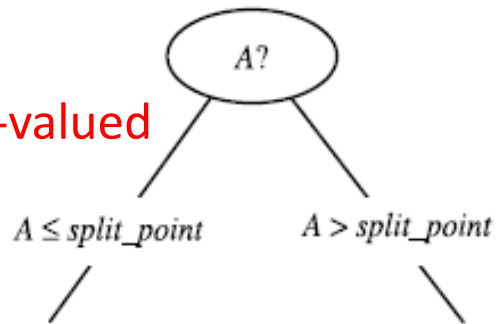
Discrete-valued

a)



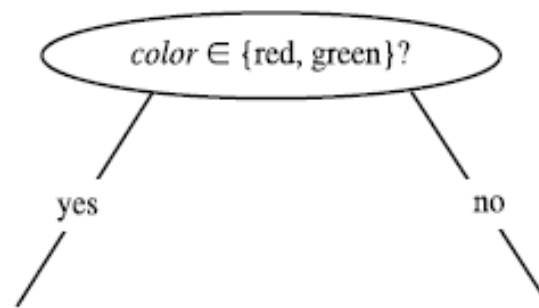
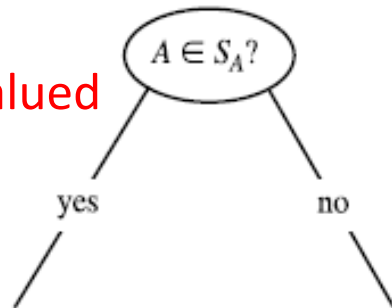
Continuous-valued

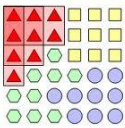
b)



Discrete-valued

c)



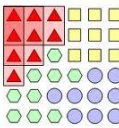


Attribute selection measures

Table 6.1 Class-labeled training tuples from the *AlIElectronics* customer database.

<i>RID</i>	<i>age</i>	<i>income</i>	<i>student</i>	<i>credit_rating</i>	<i>Class: buys_computer</i>
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

Information Gain, gain ratio and gini index



Attribute selection measures

• Information Gain

the expected information needed to classify a tuple in D : (entropy of D)

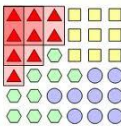
$$Info(D) = - \sum_{i=1}^m p_i \log_2(p_i),$$

➤ $p_i = |C_{i,D}| / |D|$

➤ A having v distinct values, $\{a_1, a_2, \dots, a_v\}$

➤ D_1, D_2, \dots, D_v

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j).$$

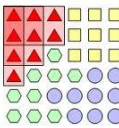


Attribute selection measures

$$Gain(A) = Info(D) - Info_A(D). \quad (\text{Choose maximum value})$$

Table 6.1 Class-labeled training tuples from the *AllElectronics* customer database.

RID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no



Example A is discrete-valued

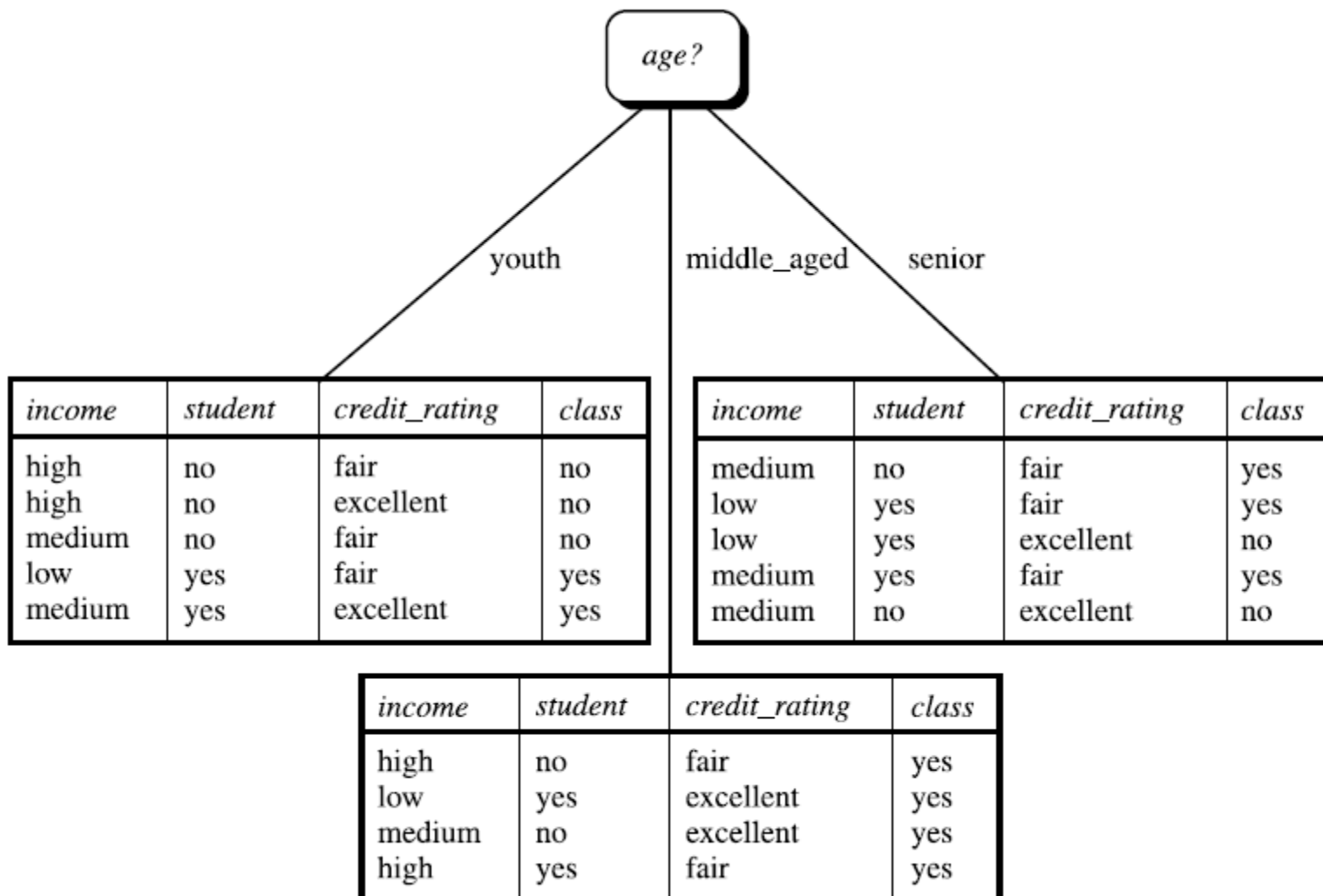
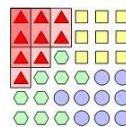
$$Info(D) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940 \text{ bits.}$$

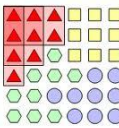
$$\begin{aligned} Info_{age}(D) &= \frac{5}{14} \times \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}\right) \\ &\quad + \frac{4}{14} \times \left(-\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4}\right) \\ &\quad + \frac{5}{14} \times \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}\right) \\ &= 0.694 \text{ bits.} \end{aligned}$$

$$Gain(age) = Info(D) - Info_{age}(D) = 0.940 - 0.694 = 0.246 \text{ bits.}$$

Gain(income)=0.029, Gain(student)=0.151,
Gain(credit_rating)=0.048

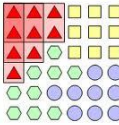
➔ Age





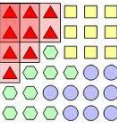
What is Prediction?

- Models continuous-valued functions, i.e. predicts unknown or missing values (numeric)
- Lost terminology of “class label attribute”, instead we use “predicted attribute”
- Viewed as a mapping or function $y = f(X)$
- Example: predict the amount (in dollars) that would be safe for the bank to loan an application



What Is Prediction?

- (Numerical) prediction is similar to classification
 - construct a model
 - use model to predict continuous or ordered value for a given input
- Prediction is different from classification
 - Classification refers to predict categorical class label
 - Prediction models continuous-valued functions
- Major method for prediction: regression
 - model the relationship between one or more *independent* or **predictor** variables and a *dependent* or **response** variable
- Regression analysis
 - Linear and multiple regression
 - Non-linear regression
 - Other regression methods: generalized linear model, Poisson regression, log-linear models, regression trees



Linear Regression

- Linear regression: involves a response variable y and a single predictor variable x

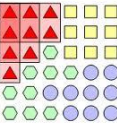
$$y = w_0 + w_1 x$$

where w_0 (y-intercept) and w_1 (slope) are regression coefficients

- Method of least squares: estimates the best-fitting straight line

$$w_1 = \frac{\sum_{i=1}^{|D|} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{|D|} (x_i - \bar{x})^2} \quad w_0 = \bar{y} - w_1 \bar{x}$$

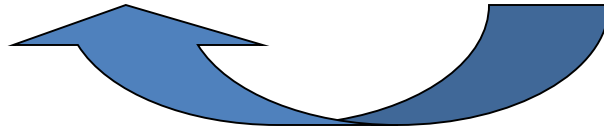
- Multiple linear regression: involves more than one predictor variable
 - Training data is of the form $(\mathbf{X}_1, y_1), (\mathbf{X}_2, y_2), \dots, (\mathbf{X}_{|D|}, y_{|D|})$
 - Ex. For 2-D data, we may have: $y = w_0 + w_1 x_1 + w_2 x_2$
 - Solvable by extension of least square method or using SAS, S-Plus
 - Many nonlinear functions can be transformed into the above



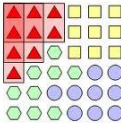
Linear Regression

A regression model is comprised of a **dependent**, or response, variable and an **independent**, or predictor, variable.

Dependent Variable = Independent Variable(s)



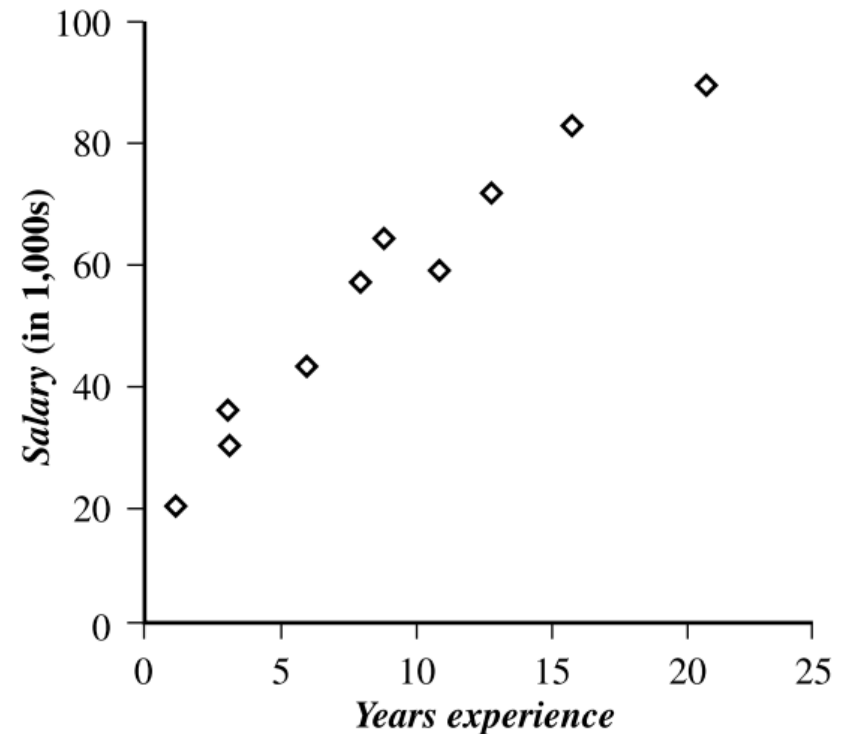
Prediction Relationship



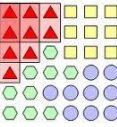
Linear Regression

Salary data.

<i>x years experience</i>	<i>y salary (in \$1000s)</i>
3	30
8	57
9	64
13	72
3	36
6	43
11	59
21	90
1	20
16	83



Plot of data: Although the points do not fall on a straight line, the overall pattern suggests a linear relationship between x (years experience) and y (salary)



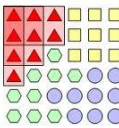
Linear Regression

Given the above data, we compute $\bar{x} = 9.1$ and $\bar{y} = 55.4$. Substituting these values into Equations (6.50) and (6.51), we get

$$w_1 = \frac{(3 - 9.1)(30 - 55.4) + (8 - 9.1)(57 - 55.4) + \dots + (16 - 9.1)(83 - 55.4)}{(3 - 9.1)^2 + (8 - 9.1)^2 + \dots + (16 - 9.1)^2} = 3.5$$

$$w_0 = 55.4 - (3.5)(9.1) = 23.6$$

Thus, the equation of the least squares line is estimated by $y = 23.6 + 3.5x$. Using this equation, we can predict that the salary of a college graduate with, say, 10 years of experience is \$58,600. ■



Nonlinear Regression

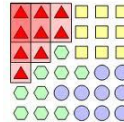
- Some nonlinear models can be modeled by a polynomial function
- A polynomial regression model can be transformed into linear regression model. For example,

$$y = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

convertible to linear with new variables: $x_2 = x^2$, $x_3 = x^3$

$$y = w_0 + w_1 x + w_2 x_2 + w_3 x_3$$

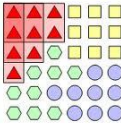
- Which is easily solved by the method of least squares using software for regression analysis
- Note that polynomial regression is a special case of multiple regression



Accuracy- Confusion Matrix

classes	buy_computer = yes	buy_computer = no	total	recognition(%)
buy_computer = yes	6954	46	7000	99.34
buy_computer = no	412	2588	3000	86.27
total	7366	2634	10000	95.42

	C_1	C_2
C_1	True positive	False negative
C_2	False positive	True negative

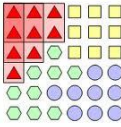


Classifier Accuracy

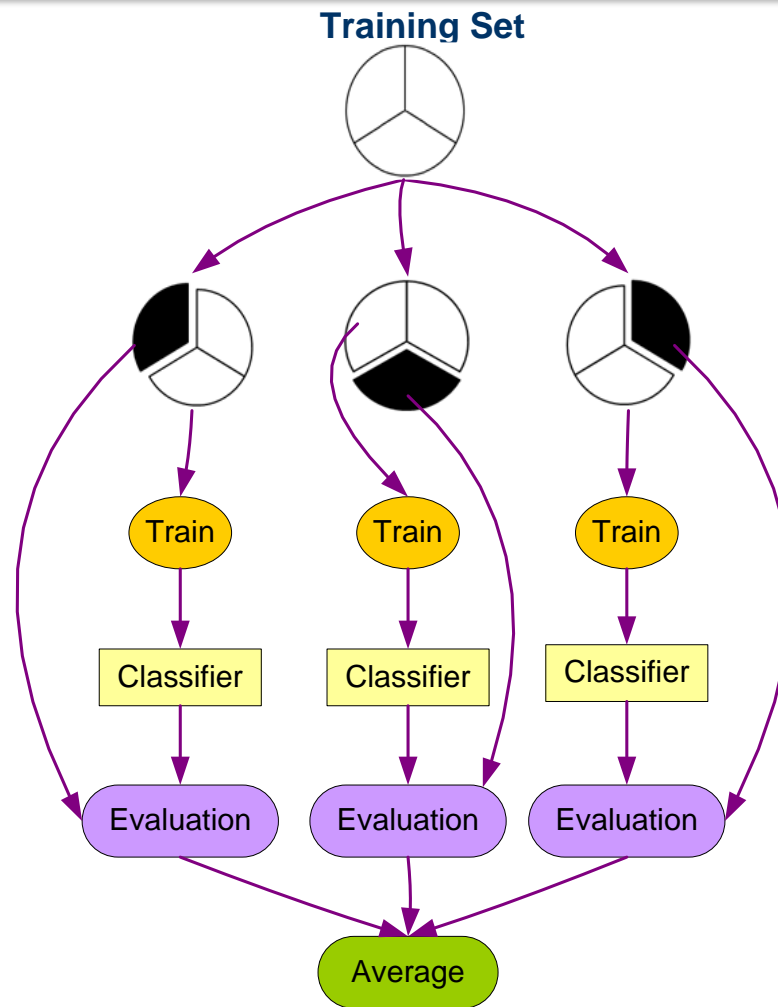
$$\text{Precision} = \frac{tp}{tp + fp}$$

$$\text{Recall} = \frac{tp}{tp + fn}$$

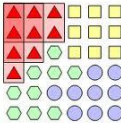
- Cross-validation (k -fold, where $k = 10$ is most popular)
 - Randomly partition the data into k *mutually exclusive* subsets, each approximately equal size
 - At i -th iteration, use D_i as test set and others as training set
 - Leave-one-out: k folds where $k = \#$ of tuples, for small sized data



Classifier Accuracy



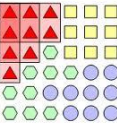
3-Crossfold Validation



Classifier Accuracy

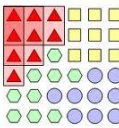
$$\text{Precision} = \frac{tp}{tp + fp}$$

$$\text{Recall} = \frac{tp}{tp + fn}$$



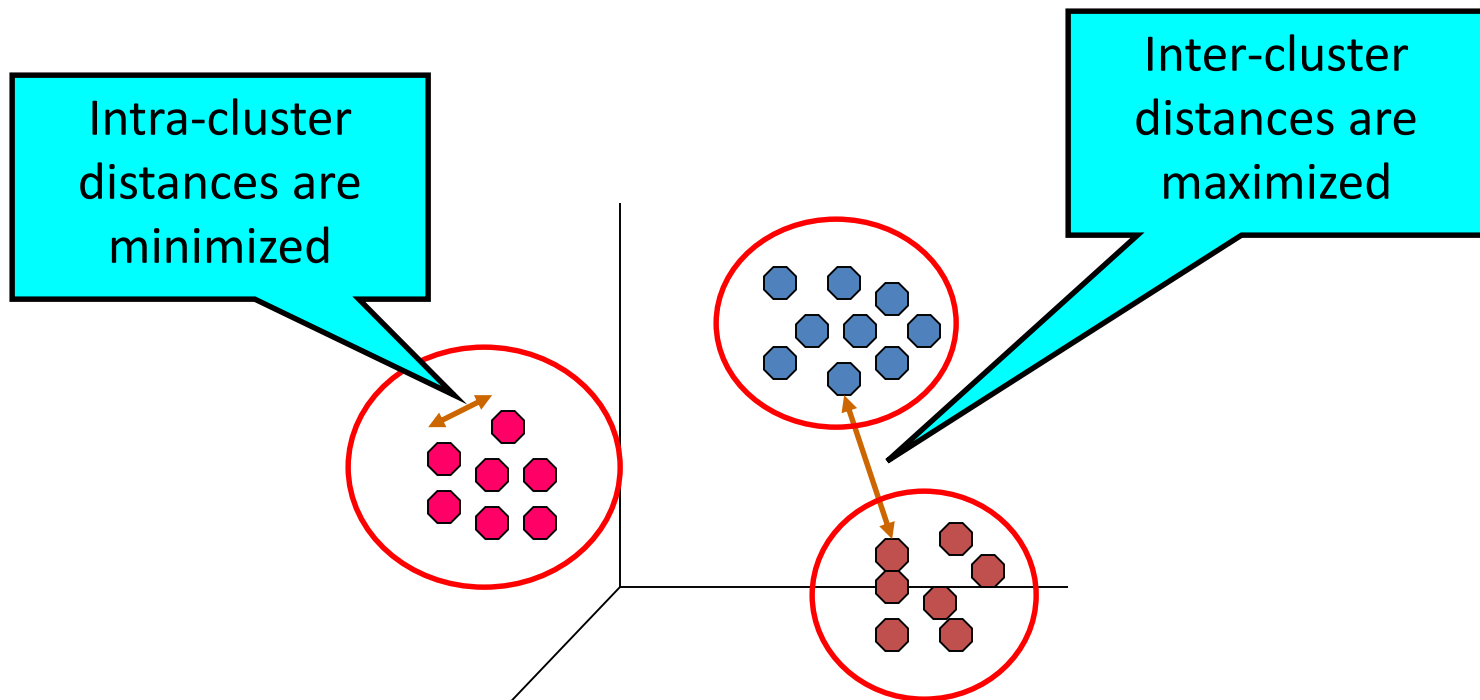
What is Cluster Analysis?

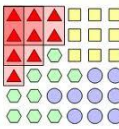
- Cluster: a collection of data objects
 - Similar to one another within the same cluster
 - Dissimilar to the objects in other clusters
- Cluster analysis
 - Grouping a set of data objects into clusters
- Clustering is **unsupervised classification**: no predefined classes
- Clustering is used:
 - As 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?

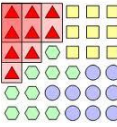
- Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups





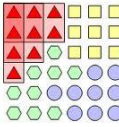
General Applications of Clustering

- 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



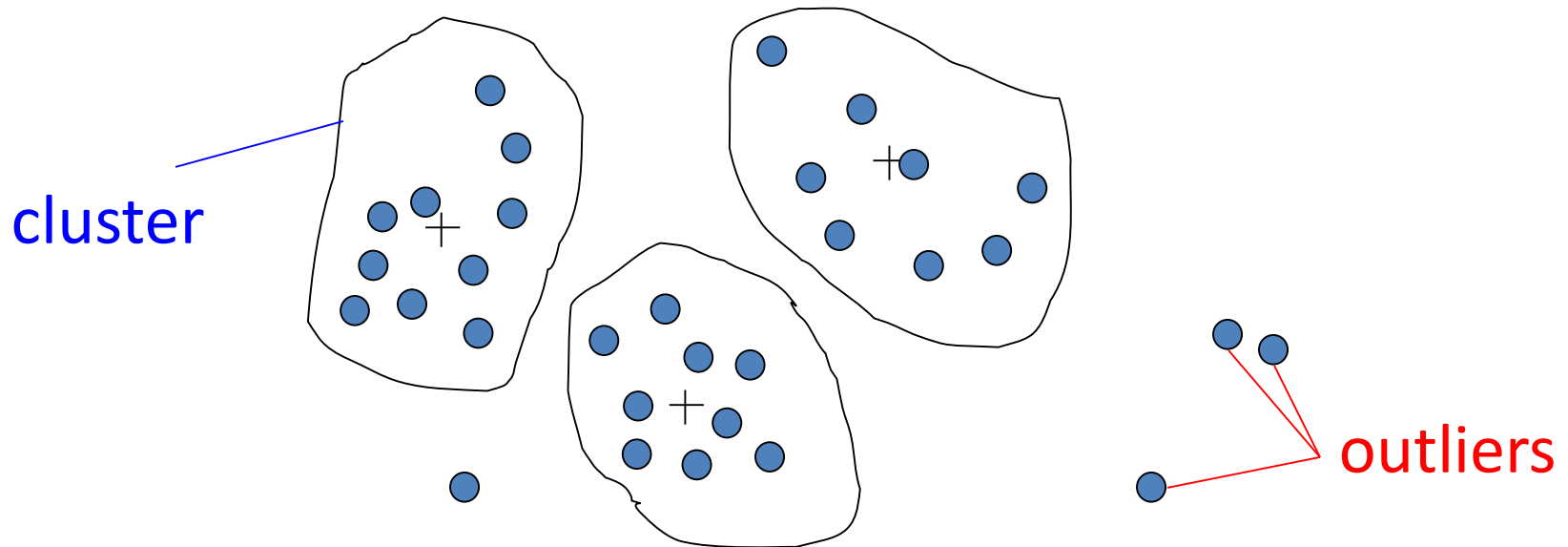
Examples of Clustering Applications

- 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

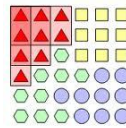


Outliers

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

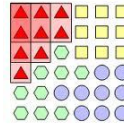


- In some applications we are interested in discovering outliers, not clusters (**outlier analysis**)



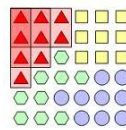
Major Clustering Approaches

- 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



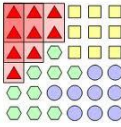
Partitioning Algorithms: Basic Concepts

- 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



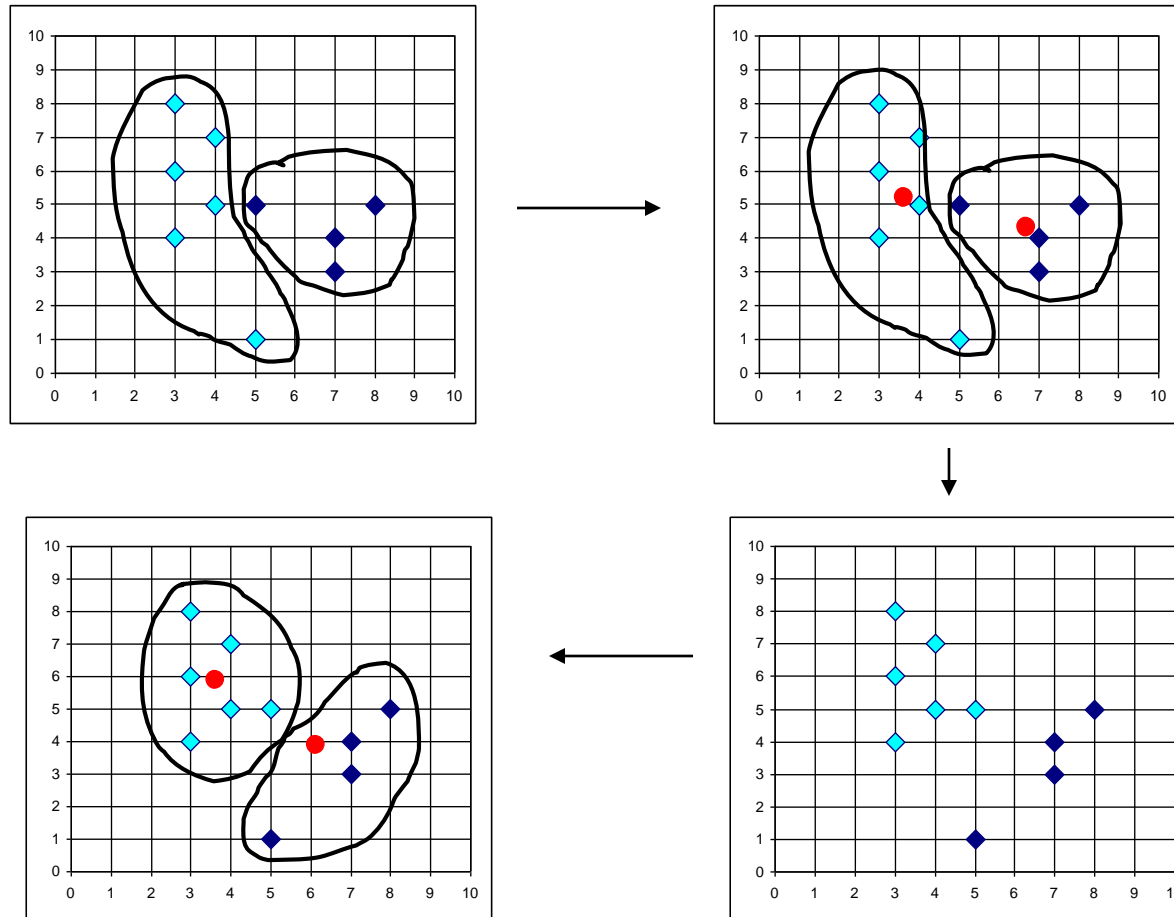
The k-means Clustering Method

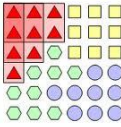
- Given k , the *k-means* algorithm is implemented in 4 steps:
 1. Partition objects into k nonempty subsets
 2. Compute seed points as the **centroids** of the clusters of the current partition. The centroid is the center (mean point) of the cluster.
 3. Assign each object to the cluster with the nearest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



The k-means Clustering Method

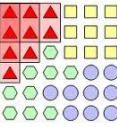
- Example





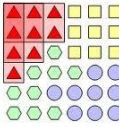
K-Means example

- 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



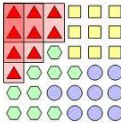
K-Means example (continued)

- 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



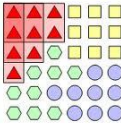
K-Means example – different starting order

- 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



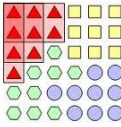
Comments on the k-means Method

- 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*



K-Medoids Method

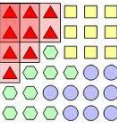
- Minimize the sensitivity of k-means to outliers
- Pick actual objects to represent clusters instead of mean values
- Each remaining object is clustered with the representative object (Medoid) to which is the most similar
- The algorithm minimizes the sum of the dissimilarities between each object and its corresponding reference point

$$E = \sum_{i=1}^k \sum_{p \in C_i} |P - O_i|$$

E: the sum of absolute error for all objects in the data set

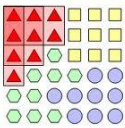
P: the data point in the space representing an object

O_i: is the representative object of cluster C_i



K-Medoids Method: The Idea

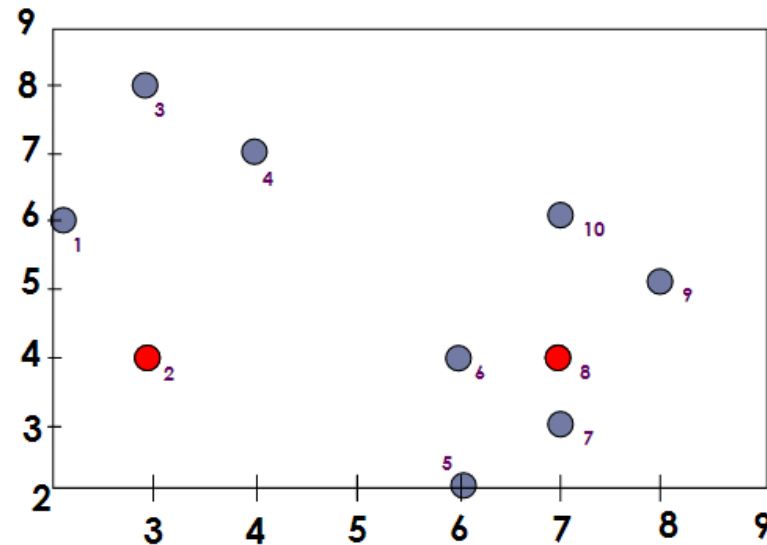
- Initial representatives are chosen randomly
- The iterative process of replacing representative objects by non-representative objects continues as long as the quality of the clustering is improved
- For each representative Object O
 - For each non-representative object R, swap O and R
- Choose the configuration with the lowest cost
- Cost function is the difference in absolute error-value if a current representative object is replaced by a non-representative object



K-Medoids Method: Example

Data Objects

	A_1	A_2
O_1	2	6
O_2	3	4
O_3	3	8
O_4	4	7
O_5	6	2
O_6	6	4
O_7	7	3
O_8	7	4
O_9	8	5
O_{10}	7	6

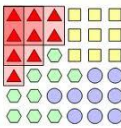


Goal: create two clusters

Choose randomly two medoids

$$O_2 = (3, 4)$$

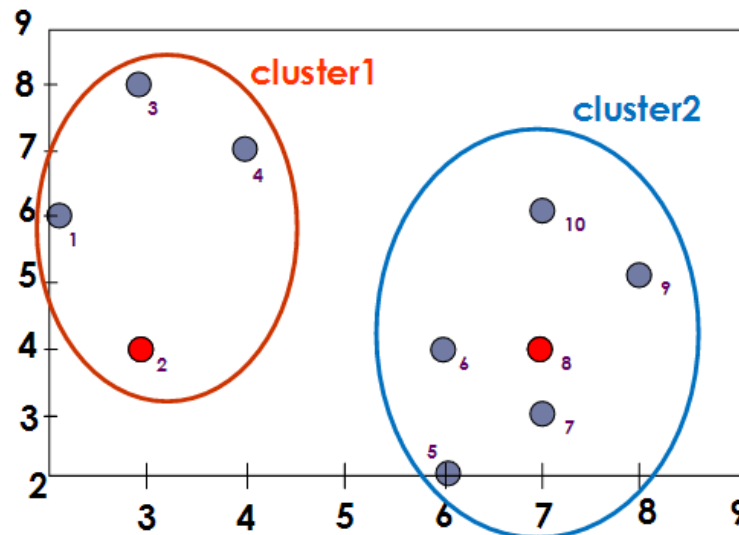
$$O_8 = (7, 4)$$



K-Medoids Method: Example

Data Objects

	A_1	A_2
O_1	2	6
O_2	3	4
O_3	3	8
O_4	4	7
O_5	6	2
O_6	6	4
O_7	7	3
O_8	7	4
O_9	8	5
O_{10}	7	6

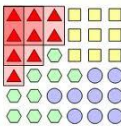


→ Assign each object to the closest representative object

→ Using L1 Metric (Manhattan), we form the following clusters

$$\text{Cluster1} = \{O_1, O_2, O_3, O_4\}$$

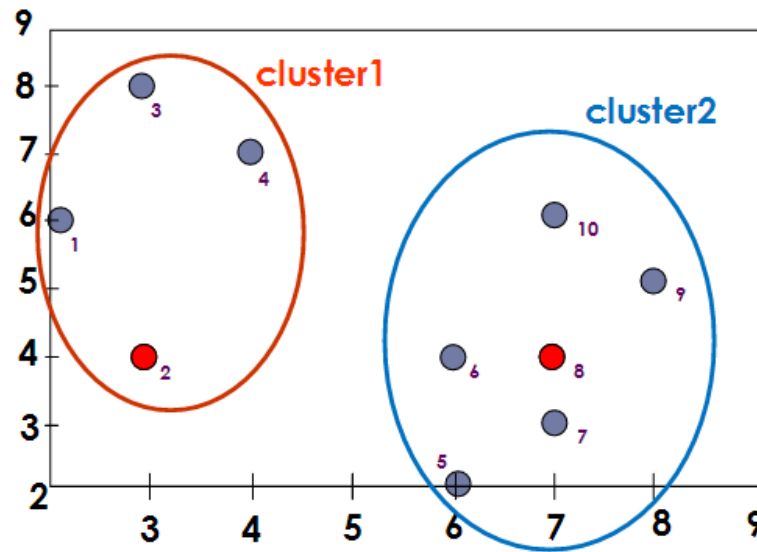
$$\text{Cluster2} = \{O_5, O_6, O_7, O_8, O_9, O_{10}\}$$



K-Medoids Method: Example

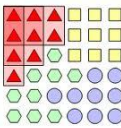
Data Objects

	A_1	A_2
O_1	2	6
O_2	3	4
O_3	3	8
O_4	4	7
O_5	6	2
O_6	6	4
O_7	7	3
O_8	7	4
O_9	8	5
O_{10}	7	6



→ Compute the absolute error criterion [for the set of Medoids (O_2, O_8)]

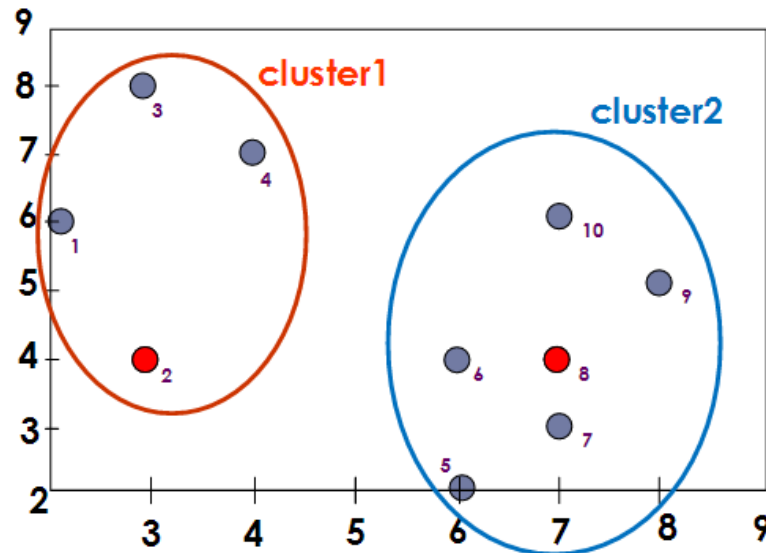
$$E = \sum_{i=1}^k \sum_{p \in C_i} |p - o_i| = |o_1 - o_2| + |o_3 - o_2| + |o_4 - o_2| + |o_5 - o_8| + |o_6 - o_8| + |o_7 - o_8| + |o_9 - o_8| + |o_{10} - o_8|$$



K-Medoids Method: Example

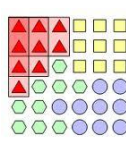
Data Objects

	A_1	A_2
O_1	2	6
O_2	3	4
O_3	3	8
O_4	4	7
O_5	6	2
O_6	6	4
O_7	7	3
O_8	7	4
O_9	8	5
O_{10}	7	6



→ The absolute error criterion [for the set of Medoids (O2,O8)]

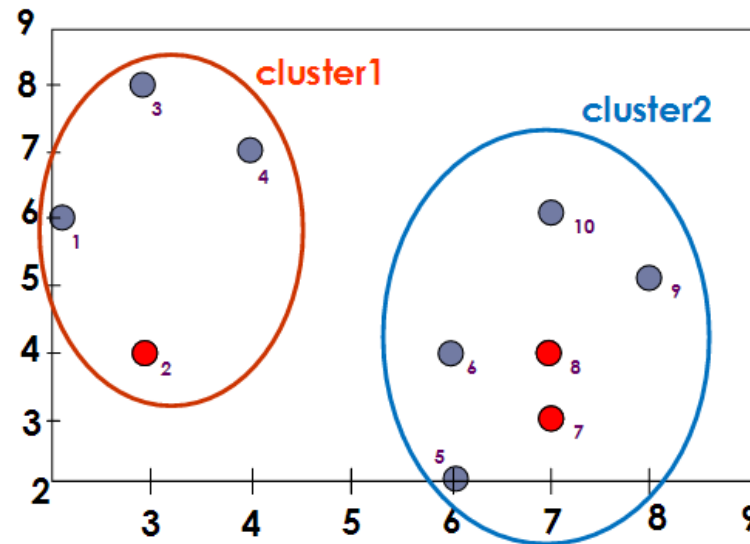
$$E = (3 + 4 + 4) + (3 + 1 + 1 + 2 + 2) = 20$$



K-Medoids Method: Example

Data Objects

	A_1	A_2
O_1	2	6
O_2	3	4
O_3	3	8
O_4	4	7
O_5	6	2
O_6	6	4
O_7	7	3
O_8	7	4
O_9	8	5
O_{10}	7	6

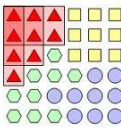


→ Choose a random object O_7

→ Swap O_8 and O_7

→ Compute the absolute error criterion [for the set of Medoids (O_2, O_7)]

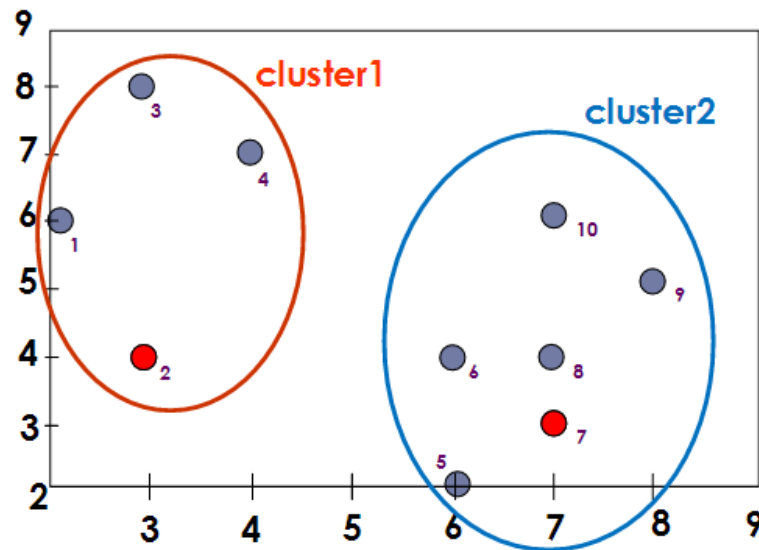
$$E = (3 + 4 + 4) + (2 + 2 + 1 + 3 + 3) = 22$$



K-Medoids Method: Example

Data Objects

	A_1	A_2
O_1	2	6
O_2	3	4
O_3	3	8
O_4	4	7
O_5	6	2
O_6	6	4
O_7	7	3
O_8	7	4
O_9	8	5
O_{10}	7	6

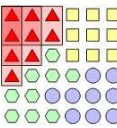


→ Compute the cost function

Absolute error [for O_2, O_7] – Absolute error [O_2, O_8]

$$S = 22 - 20$$

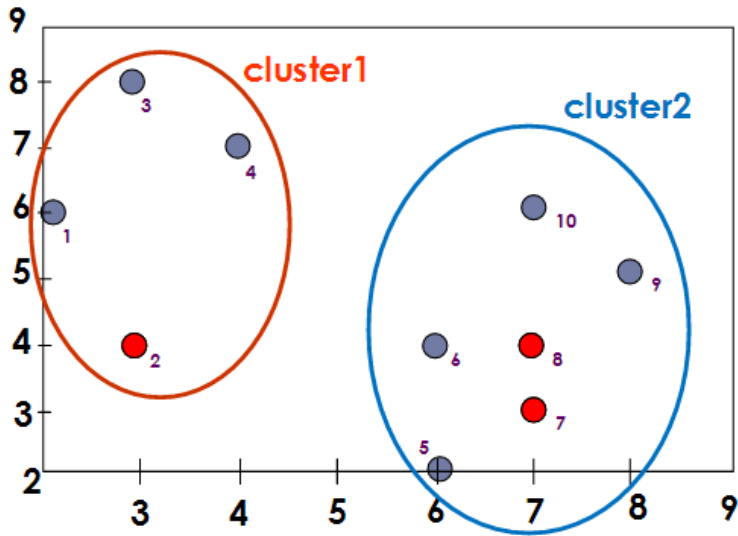
$S > 0 \Rightarrow$ it is a bad idea to replace O_8 by O_7



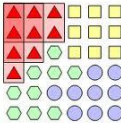
K-Medoids Method: Example

Data Objects

	A_1	A_2
O_1	2	6
O_2	3	4
O_3	3	8
O_4	4	7
O_5	6	2
O_6	6	4
O_7	7	3
O_8	7	4
O_9	8	5
O_{10}	7	6

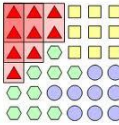


- ▶ In this example, changing the medoid of cluster 2 did not change the assignments of objects to clusters.
- ▶ What are the possible cases when we replace a medoid by another object?

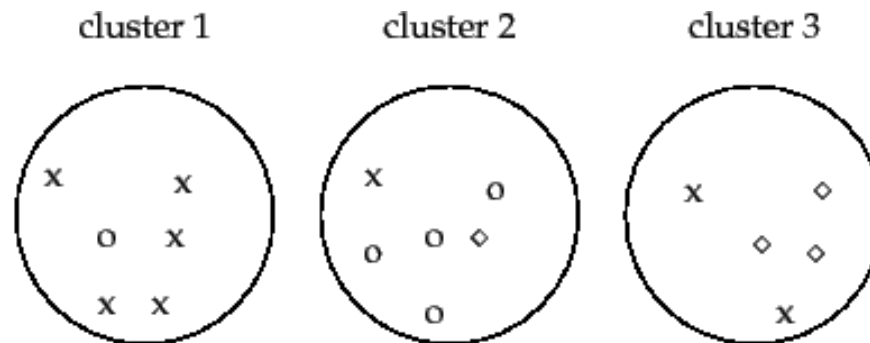


Validation Criteria : Clustering

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



Clustering: Purity

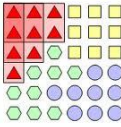


► **Figure 16.1** Purity as an external evaluation criterion for cluster quality. Majority class and number of members of the majority class for the three clusters are: x, 5 (cluster 1); o, 4 (cluster 2); and \diamond , 3 (cluster 3). Purity is $(1/17) \times (5 + 4 + 3) \approx 0.71$.

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

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

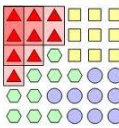
A true positive (TP) decision assigns two similar documents to the same cluster, a true negative (TN) decision assigns two dissimilar documents to different clusters. There are two types of errors we can commit. A (FP) decision assigns two dissimilar documents to the same cluster. A (FN) decision assigns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

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

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

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

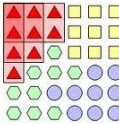
$$FP = 40 - 20 = 20$$



Clustering: Rand Index

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

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



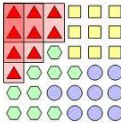
Clustering: F-Measure

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

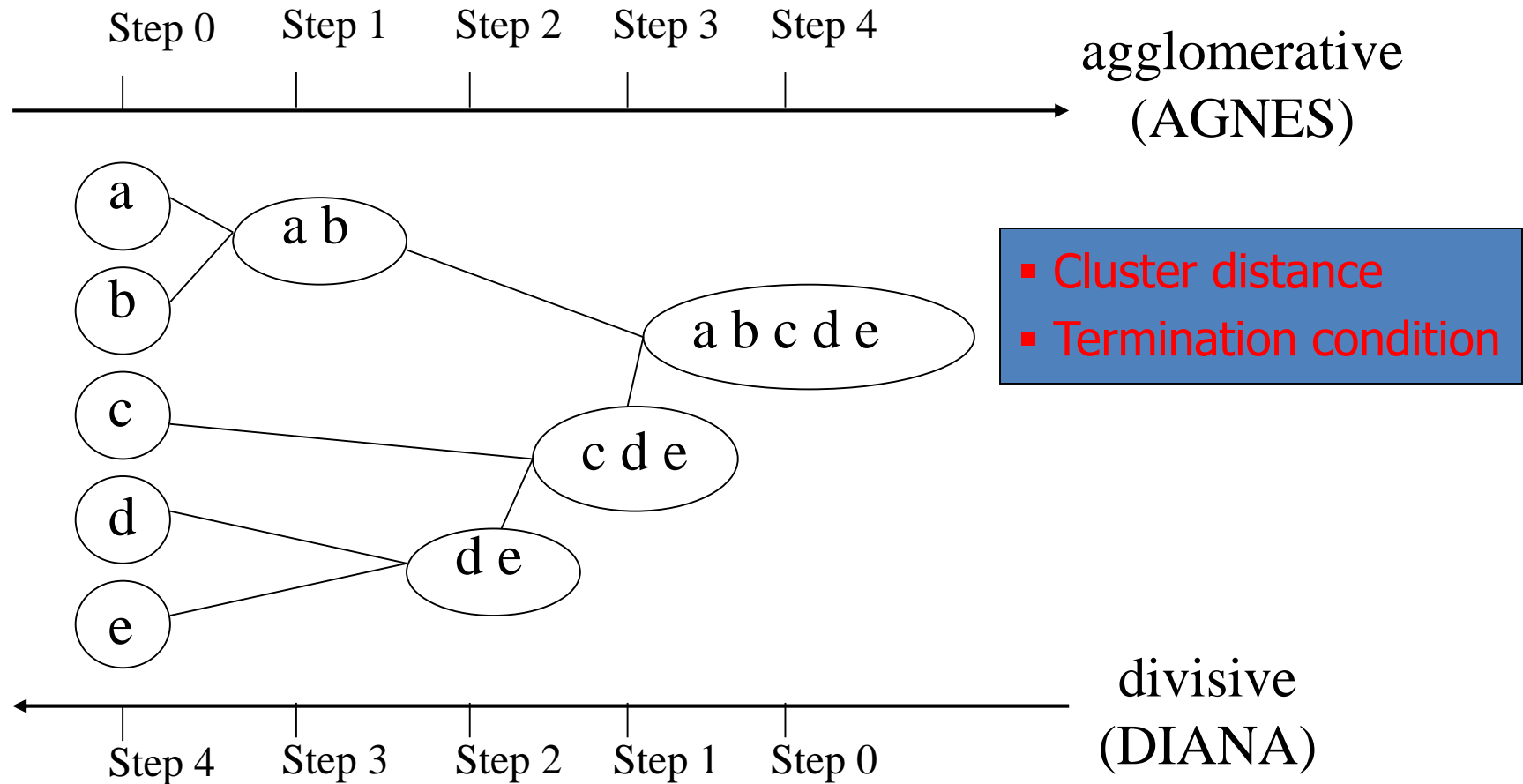
$$P = \frac{TP}{TP + FP} \quad R = \frac{TP}{TP + FN} \quad F_{\beta} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

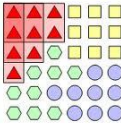
$$P = 20/40 = 0.5 \quad R = 20/44 \approx 0.455 \quad \underline{F_1 \approx 0.48} \quad \beta = 1$$
$$\underline{F_5 \approx 0.456} \quad \beta = 5$$

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



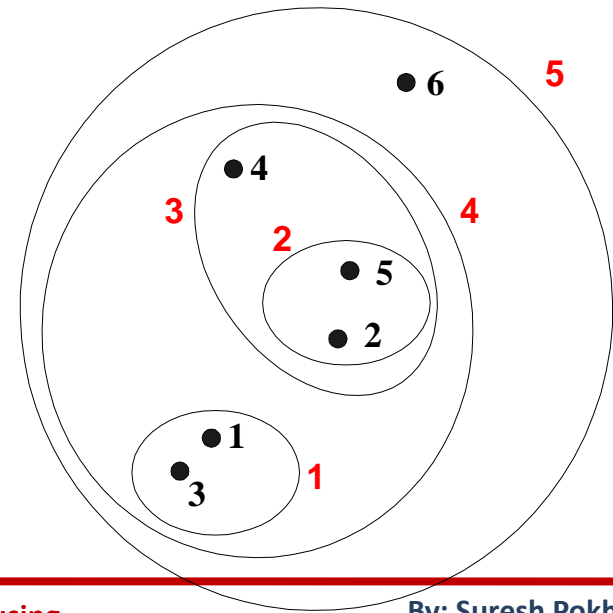
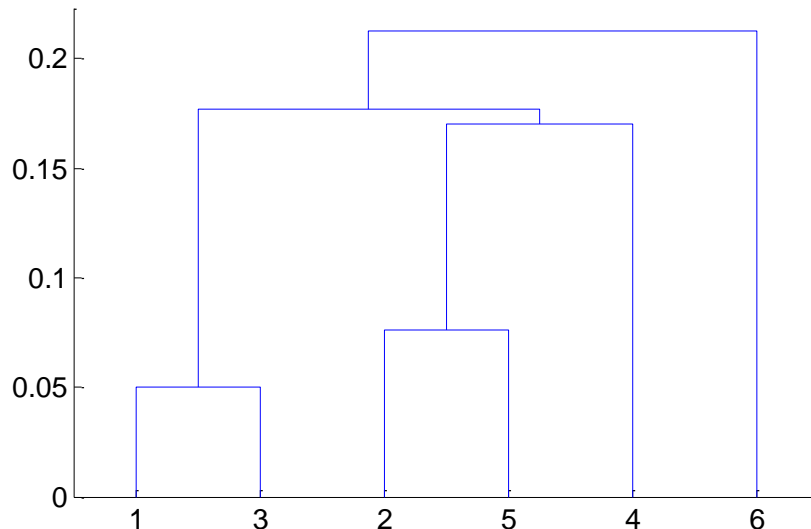
Hierarchical Clustering

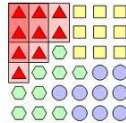




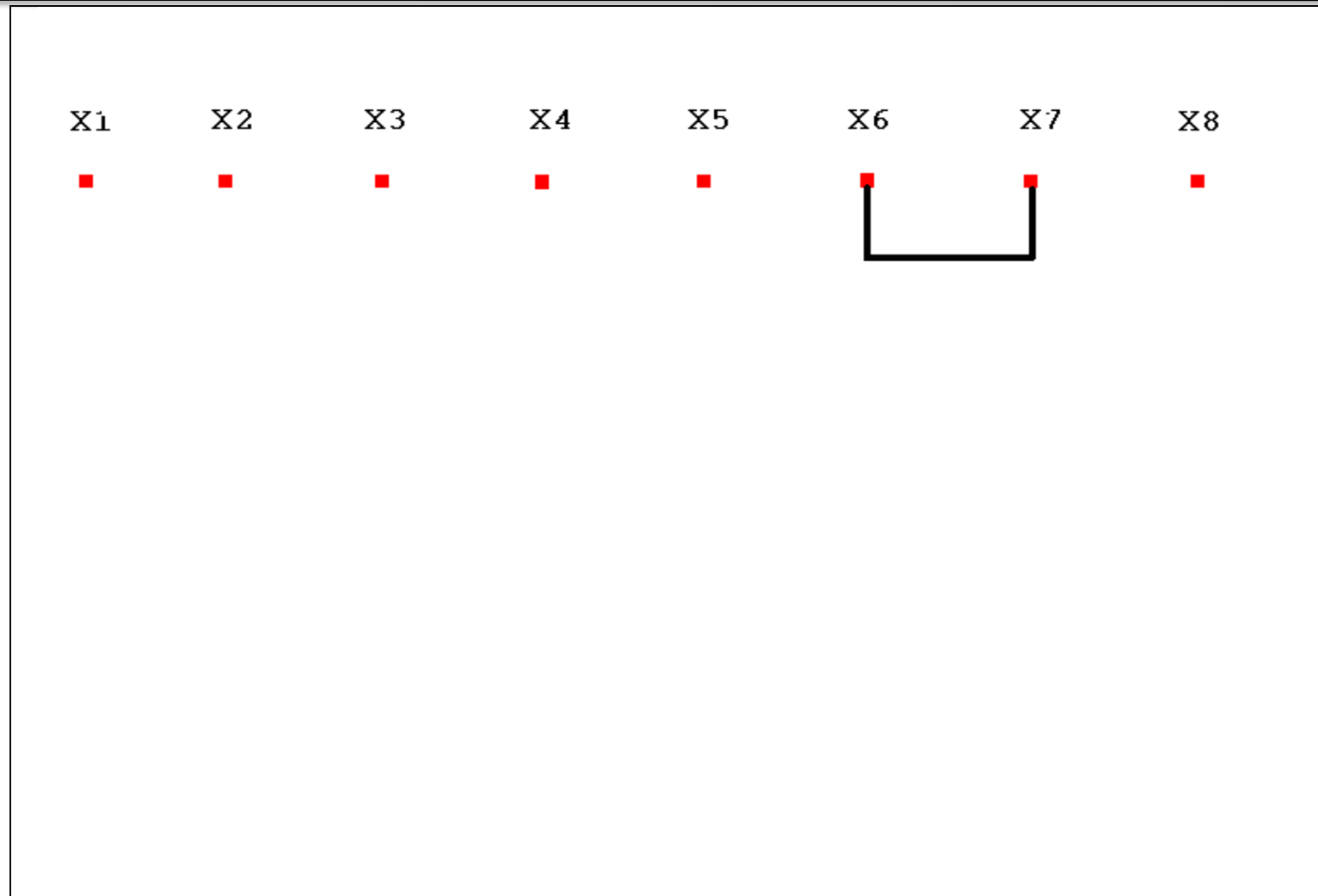
Hierarchical Clustering

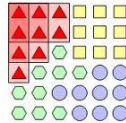
- Produces a set of nested clusters organized as a hierarchical tree
- Can be visualized as a “dendrogram”
 - A tree like diagram that records the sequences of merges or splits



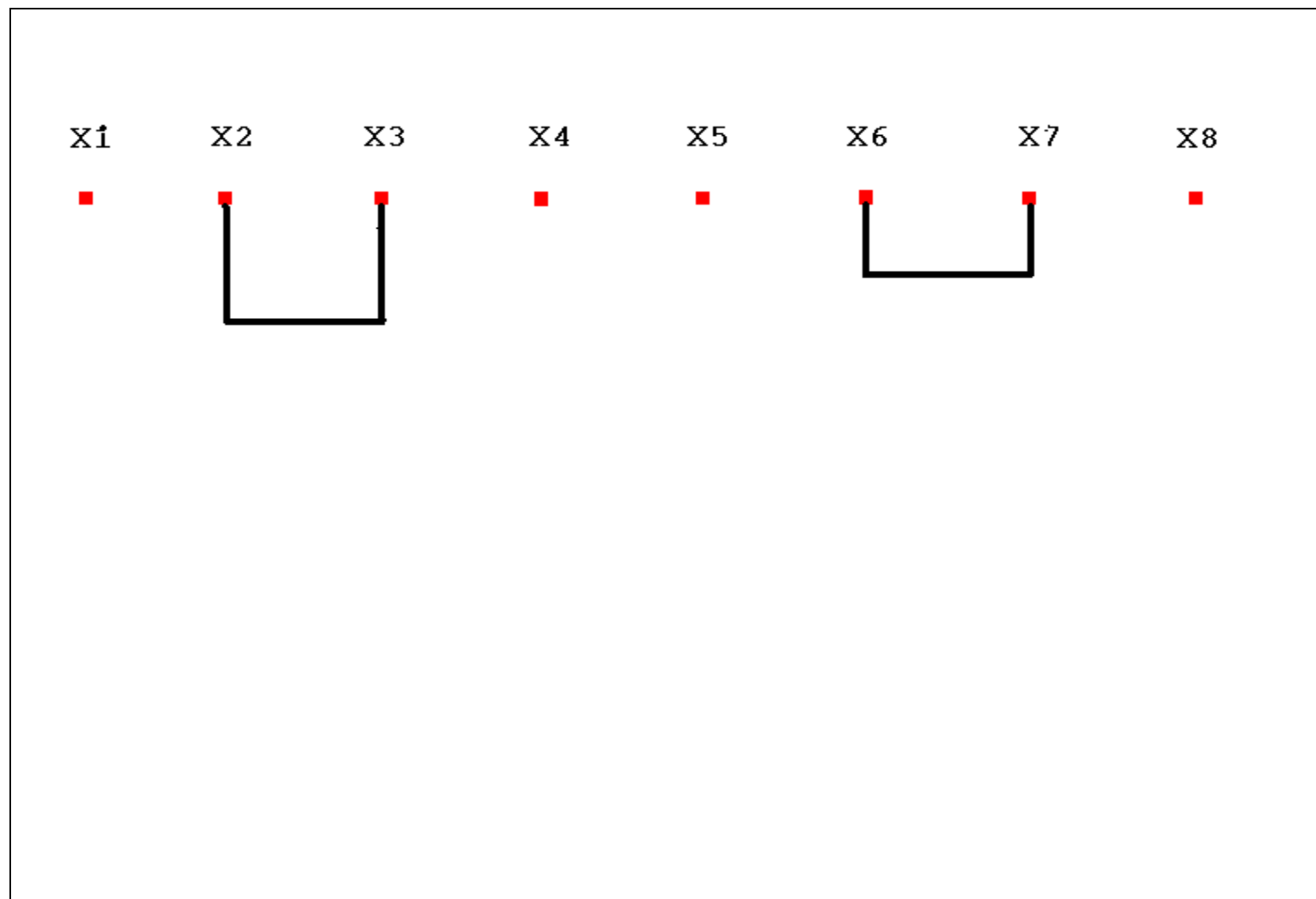


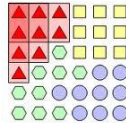
Nearest Neighbor, Level 2, $k = 7$ clusters.



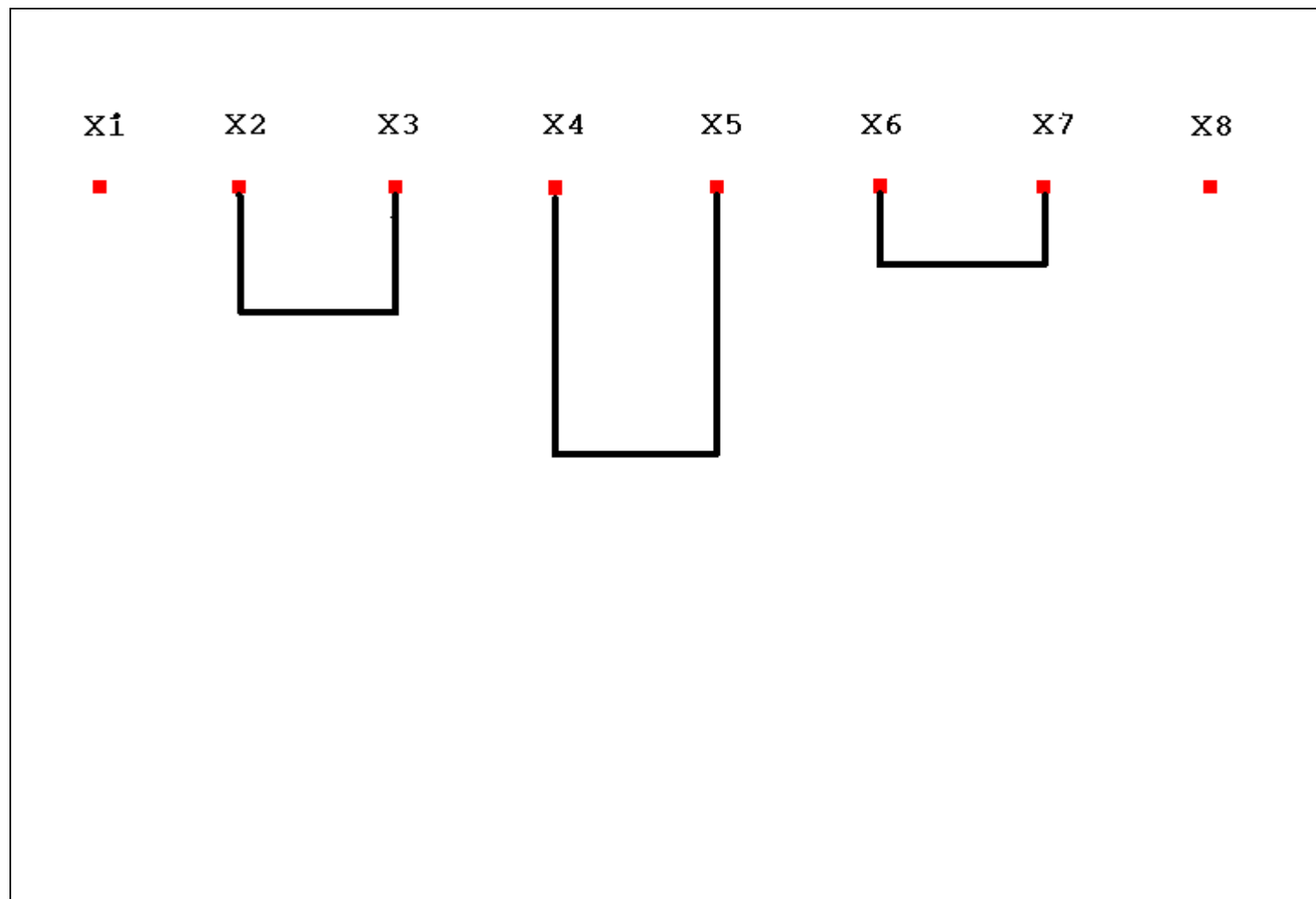


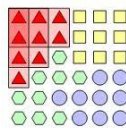
Nearest Neighbor, Level 3, $k = 6$ clusters.



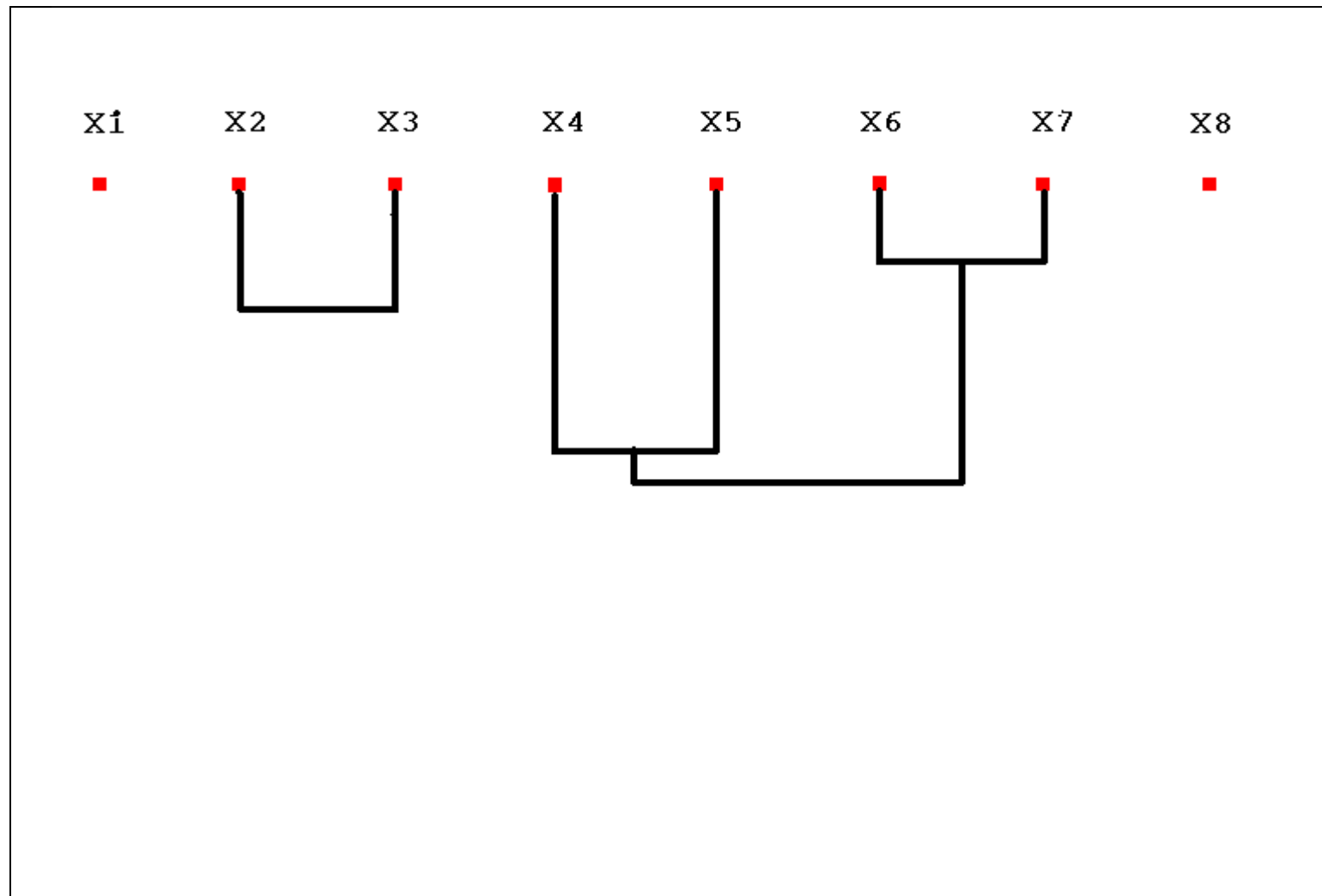


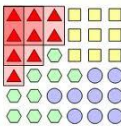
Nearest Neighbor, Level 4, $k = 5$ clusters.



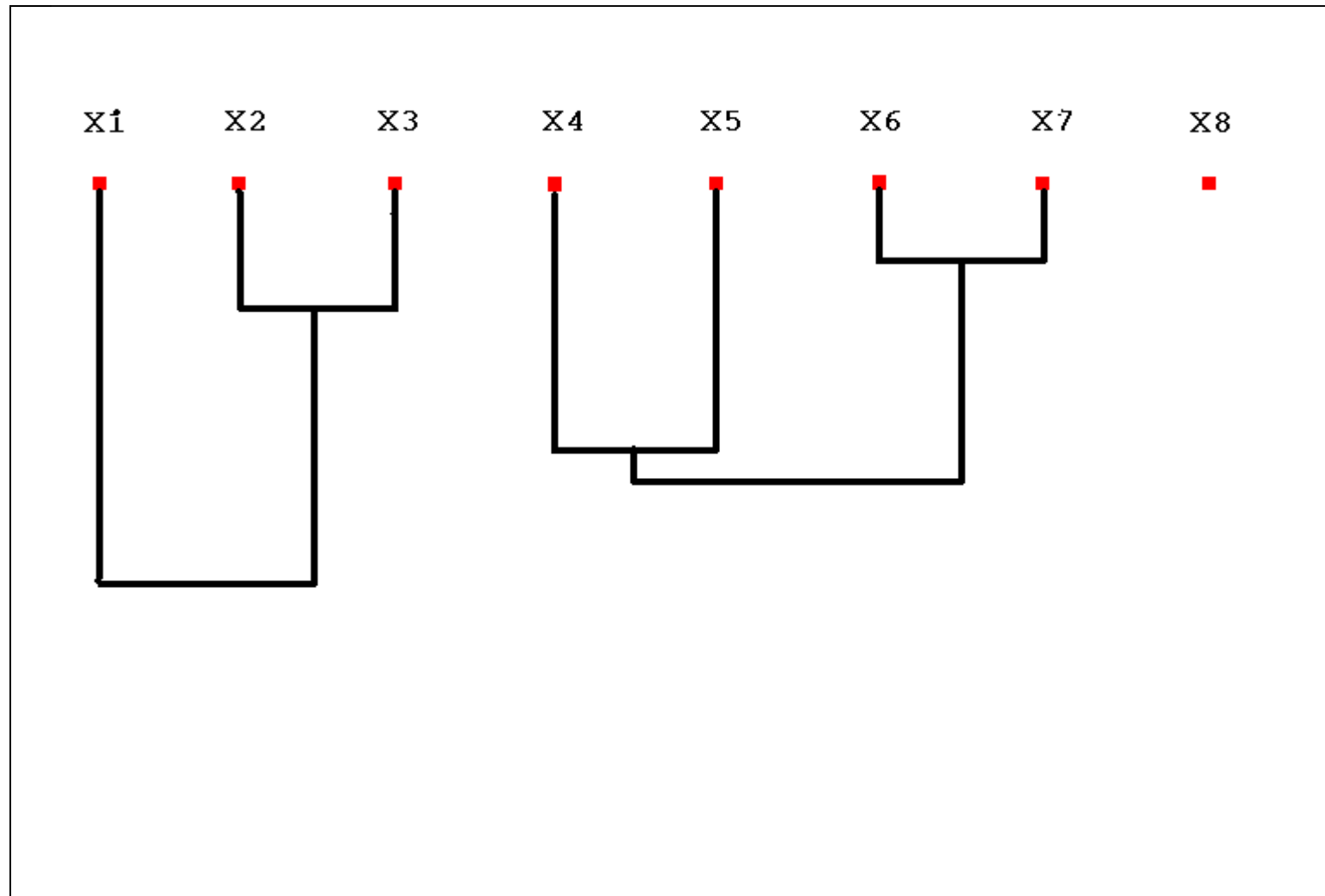


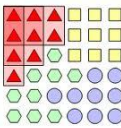
Nearest Neighbor, Level 5, $k = 4$ clusters.



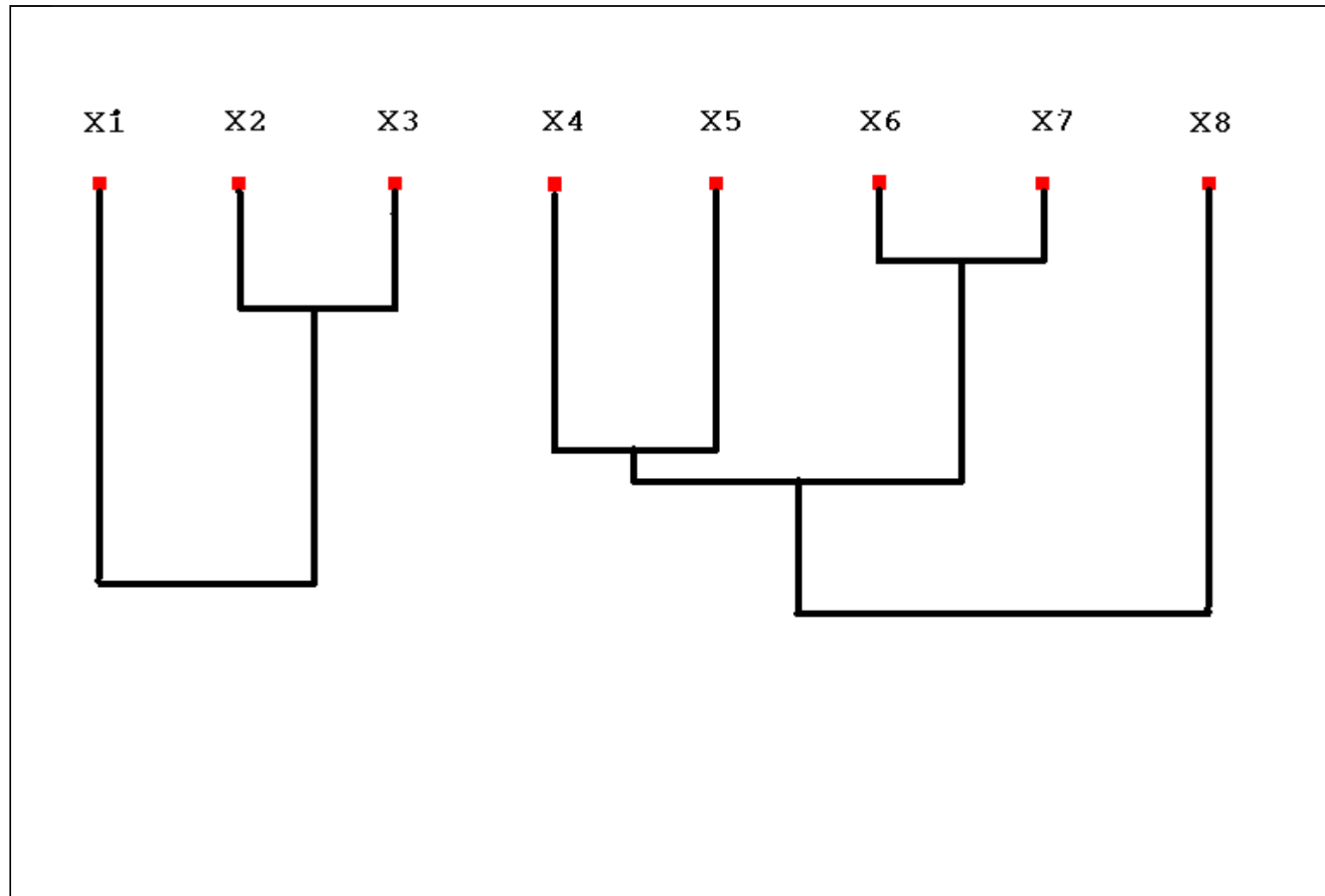


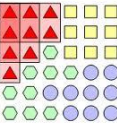
Nearest Neighbor, Level 6, $k = 3$ clusters.



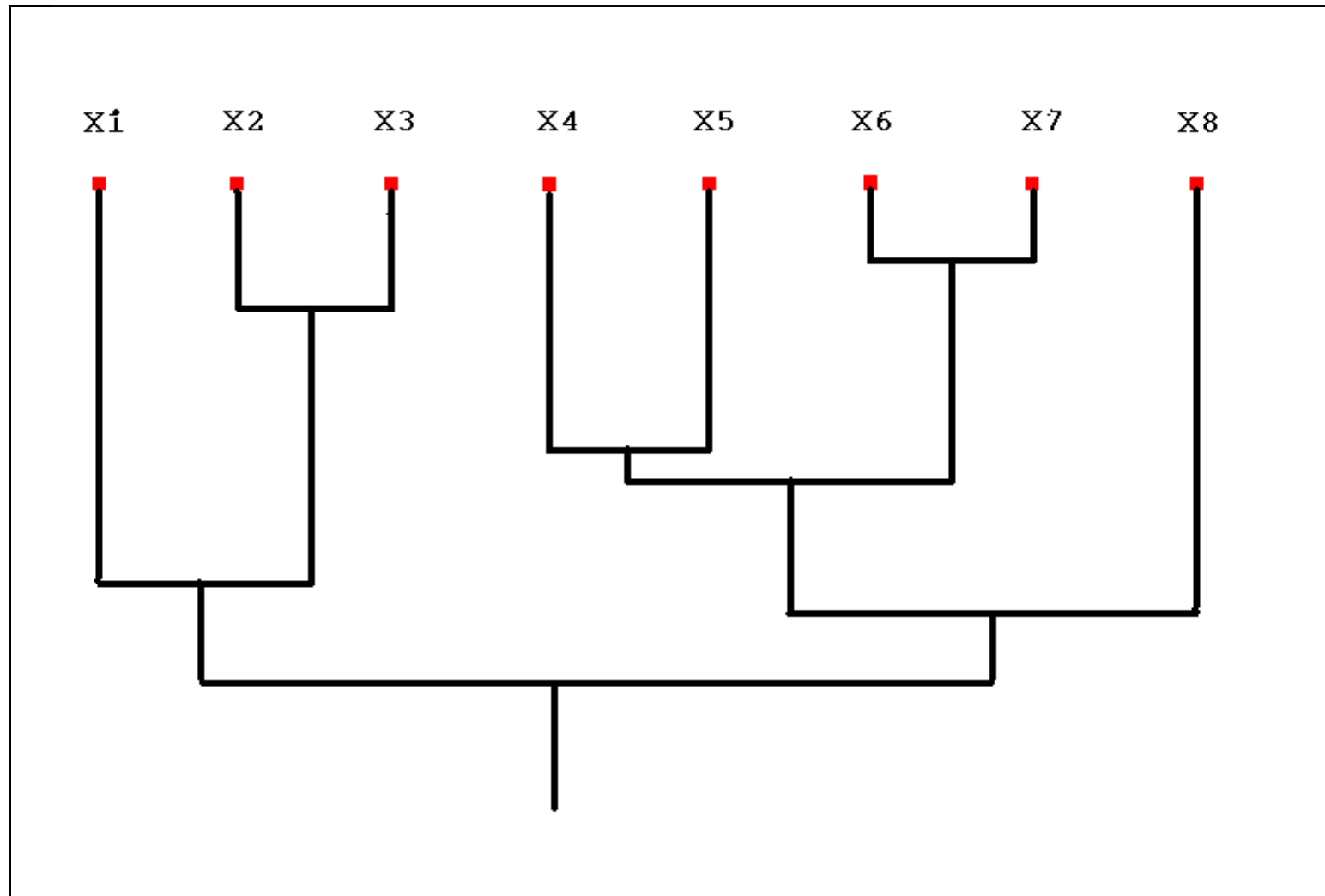


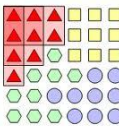
Nearest Neighbor, Level 7, $k = 2$ clusters.





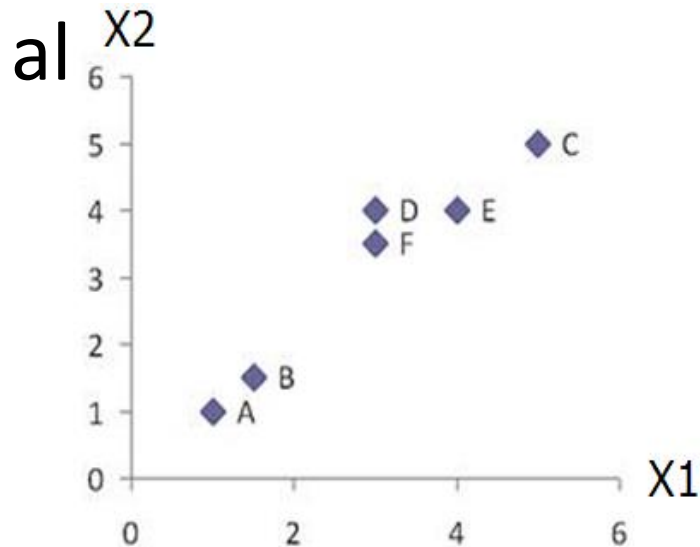
Nearest Neighbor, Level 8, k = 1 cluster.





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

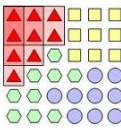
$$d_{AB} = \left((1-1.5)^2 + (1-1.5)^2 \right)^{\frac{1}{2}} = \sqrt{\frac{1}{2}} = 0.7071$$

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

Euclidean distance

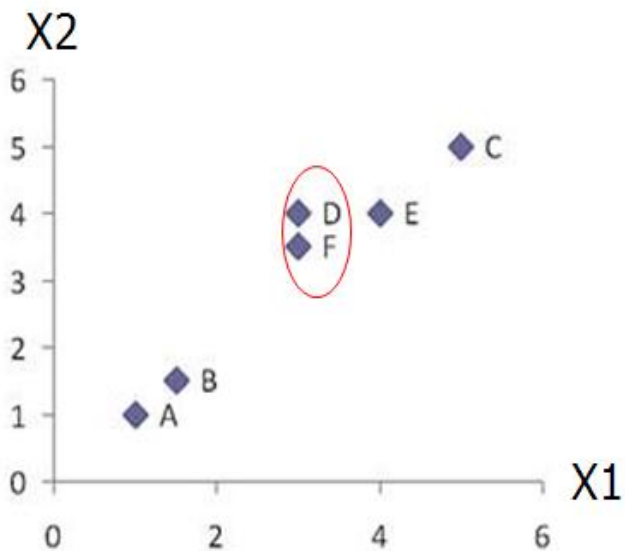
Dist	A	B	C	D	E	F
A	0.00	0.71	5.66	3.61	4.24	3.20
B	0.71	0.00	4.95	2.92	3.54	2.50
C	5.66	4.95	0.00	2.24	1.41	2.50
D	3.61	2.92	2.24	0.00	1.00	0.50
E	4.24	3.54	1.41	1.00	0.00	1.12
F	3.20	2.50	2.50	0.50	1.12	0.00

distance matrix



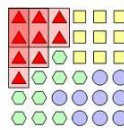
Example and Demo

- 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

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

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

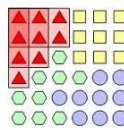
$$d_{(D,F) \rightarrow C} = \min(d_{DC}, d_{FC}) = \min(2.24, 2.50) = 2.24$$

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

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

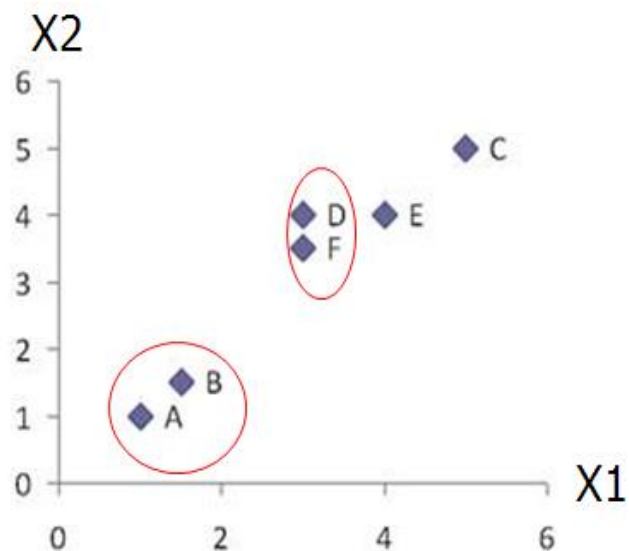
Min Distance (Single Linkage)

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

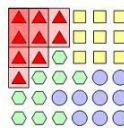
- Merge two closest clusters



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

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

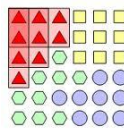
$$d_{(D,F) \rightarrow (A,B)} = \min(d_{DA}, d_{DB}, d_{FA}, d_{FB}) = \min(3.61, 2.92, 3.20, 2.50) = 2.50$$

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

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

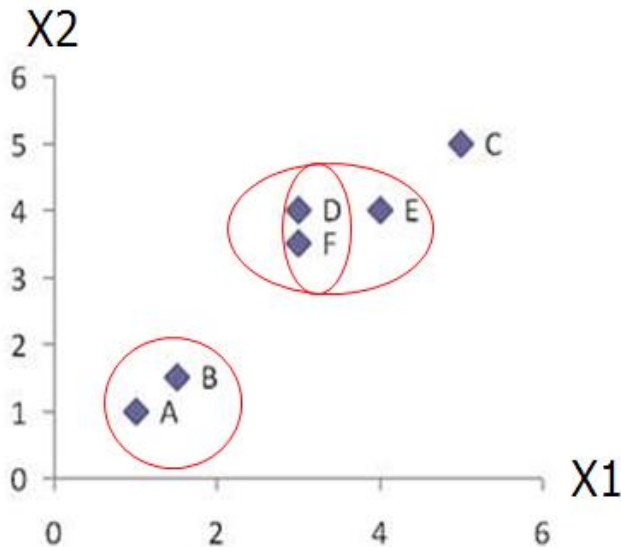
Min Distance (Single Linkage)

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

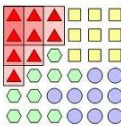


Min Distance (Single Linkage)

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

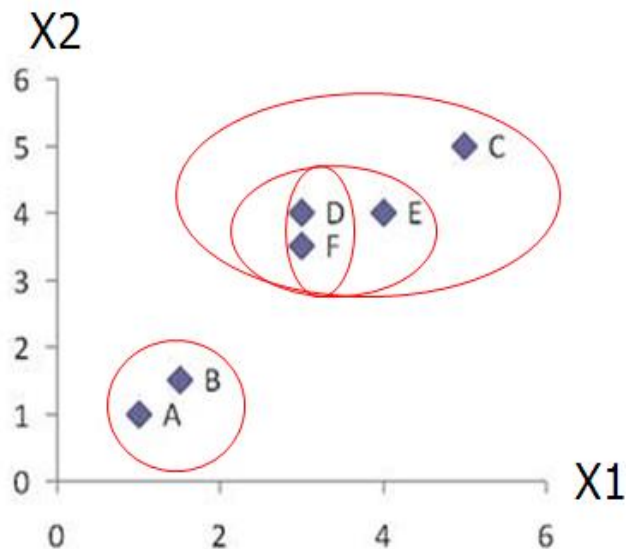
Min Distance (Single Linkage)

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

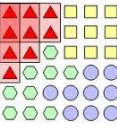


Min Distance (Single Linkage)

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)

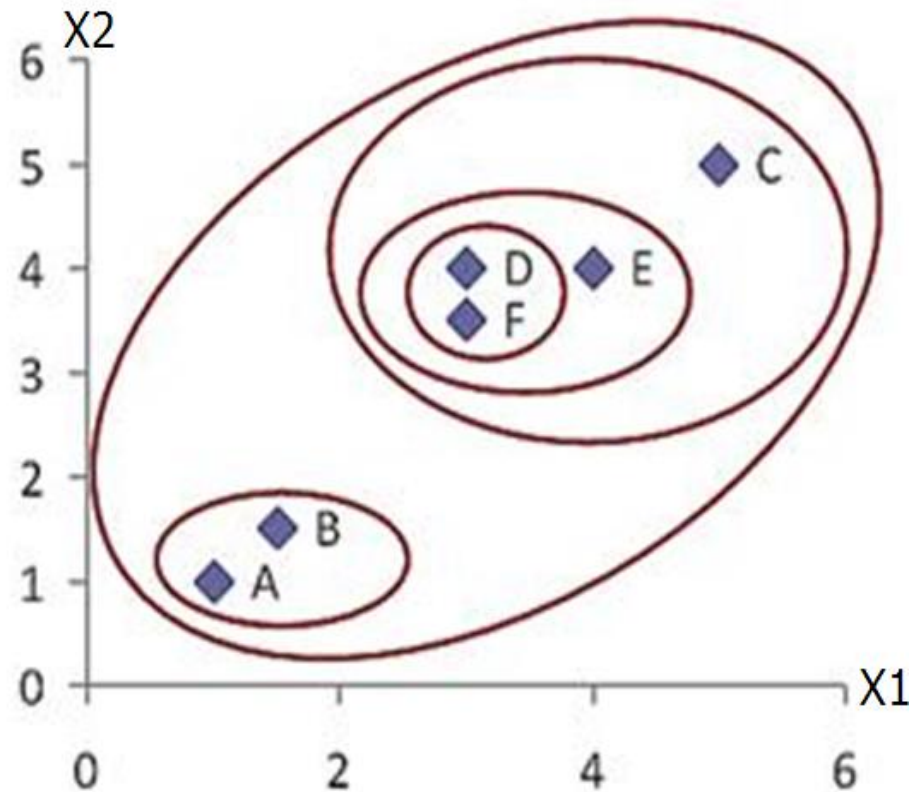
Dist	(A,B)	((D, F), E), C
(A,B)	0.00	2.50
((D, F), E), C	2.50	0.00

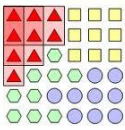


Example and Demo

- Final result (meeting termination condition)

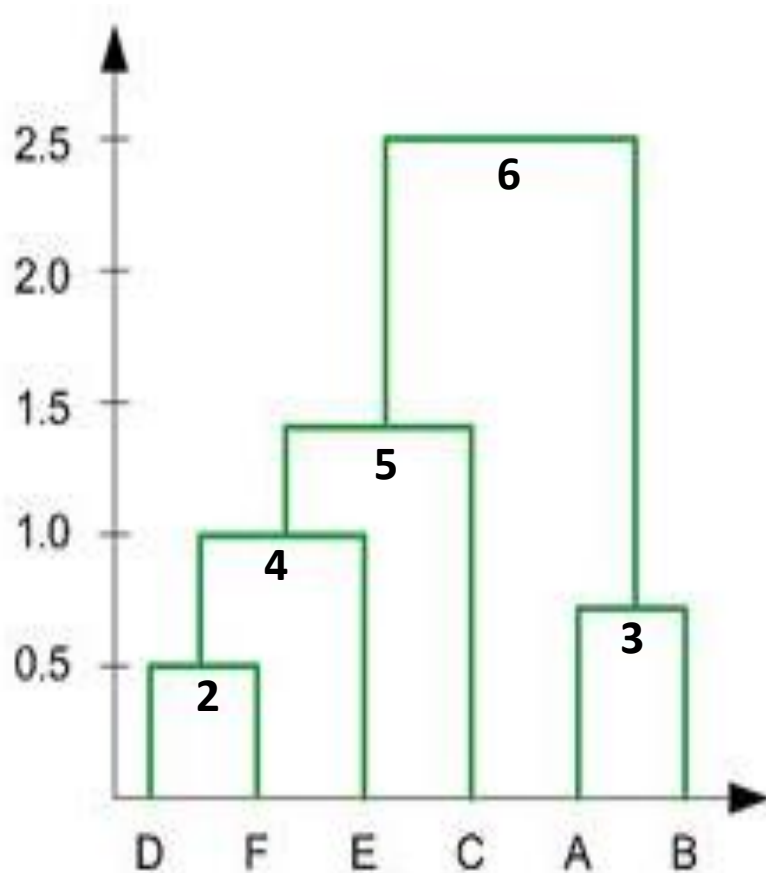
	X1	X2
A	1	1
B	1.5	1.5
C	5	5
D	3	4
E	4	4
F	3	3.5



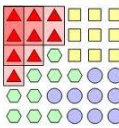


Example and Demo

- **Dendrogram tree** representation

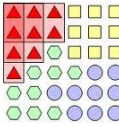


1. In the beginning we have 6 clusters: A, B, C, D, E and F
2. We merge cluster D and F into cluster (D, F) at distance 0.50
3. We merge cluster A and cluster B into (A, B) at distance 0.71
4. We merge cluster E and (D, F) into ((D, F), E) at distance 1.00
5. We merge cluster ((D, F), E) and C into (((D, 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



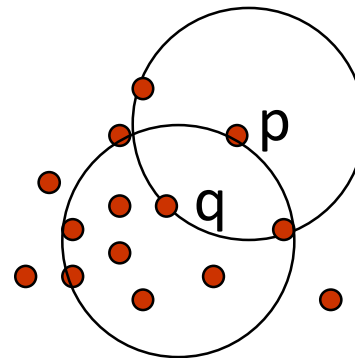
Density-Based Clustering Methods

- 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)

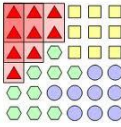


Density-Based Clustering: Background

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



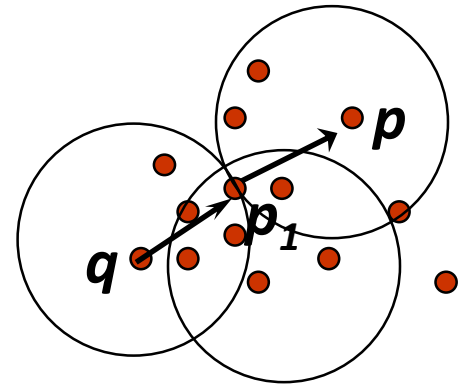
MinPts = 5
Eps = 1 cm



Density-Based Clustering: Background (II)

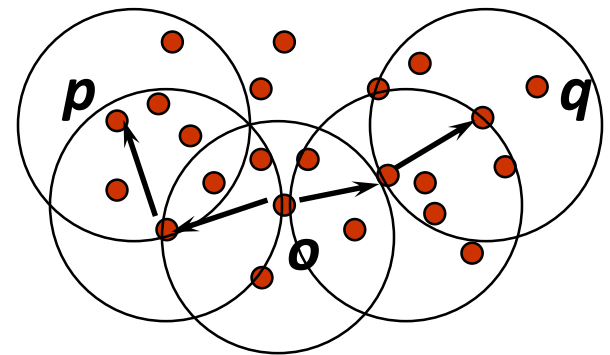
- **Density-reachable:**

- 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**

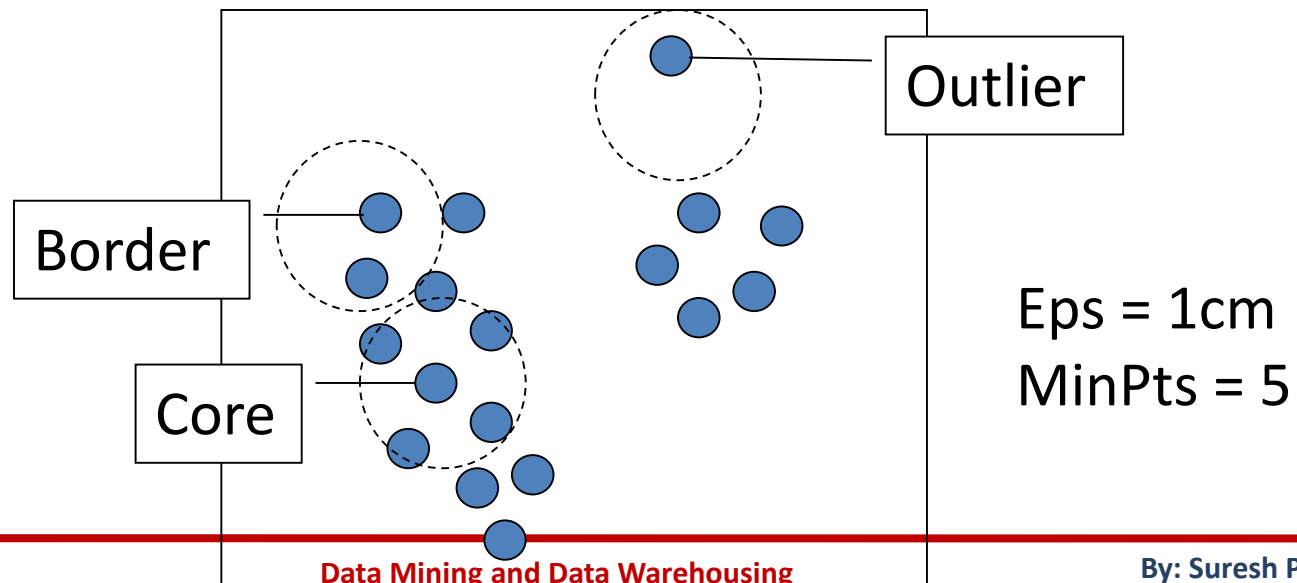
- A point p is density-connected to a point q wrt. Eps , $MinPts$ if there is a point o such that both, p and q are density-reachable from o wrt. Eps and $MinPts$.

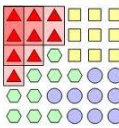




DBSCAN: Density Based Spatial Clustering of Applications with Noise

- 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.

