

(Computer Engineering and Technology) (SY B.Tech-AI-DS)

Disclaimer:

1. Information included in this slides came from multiple sources. We have tried our best to cite the sources. Please refer to the References to learn about the sources, when applicable.
2. The slides should be used only for academic purposes (e.g., in teaching a class), and should not be used for commercial purposes.

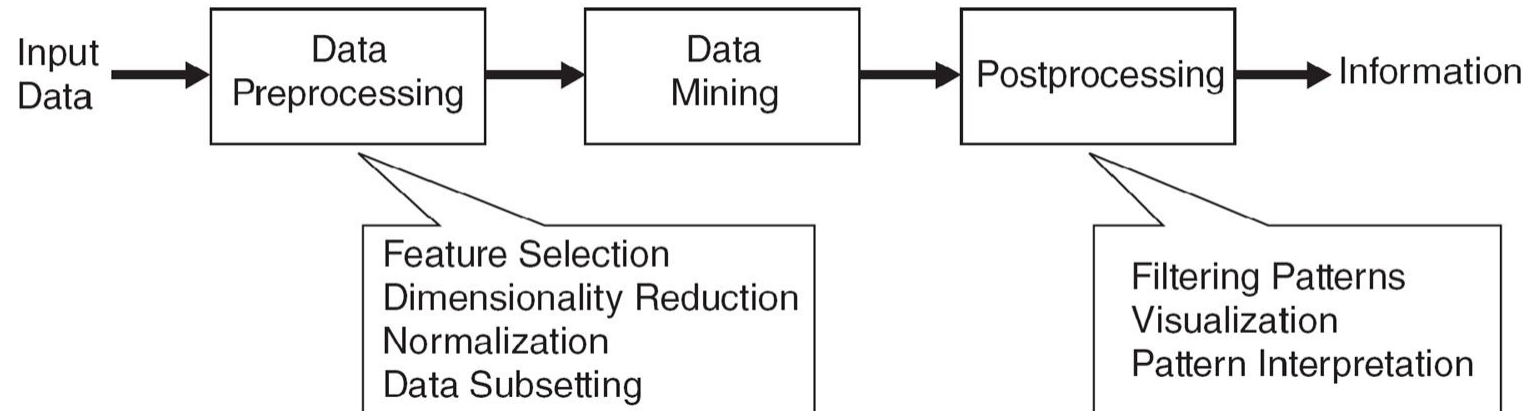
UNIT V-Syllabus

Data Mining : Classification And Clustering

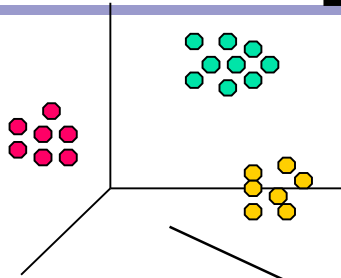
- Introduction to Data Mining, Data Mining Techniques, Supervised, Semi-Supervised, and Unsupervised Methods, Classification: Basic
- Concepts, Decision Tree Induction, Bayesian Classification
- Clustering Techniques: Basic concepts, Partition based Clustering: k-Means

Data Mining

- Data Mining is Non-trivial **extraction of implicit, previously unknown and potentially useful information from data**
- **Data Mining is Exploration & analysis**, by automatic or semi-automatic means, of large quantities of data in order to discover meaningful patterns



Data Mining Task..



Clustering

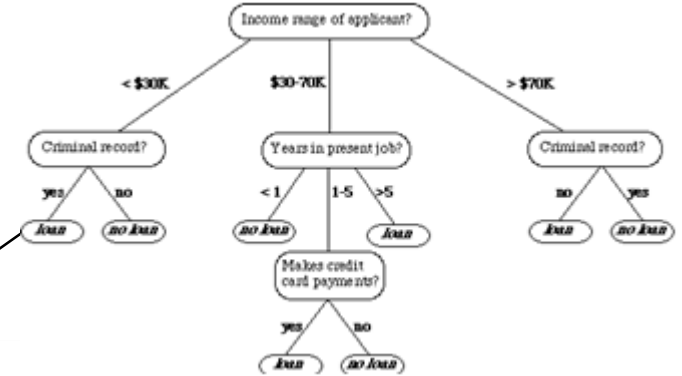
Data

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
11	No	Married	60K	No
12	Yes	Divorced	220K	No
13	No	Single	85K	Yes
14	No	Married	75K	No
15	No	Single	90K	Yes

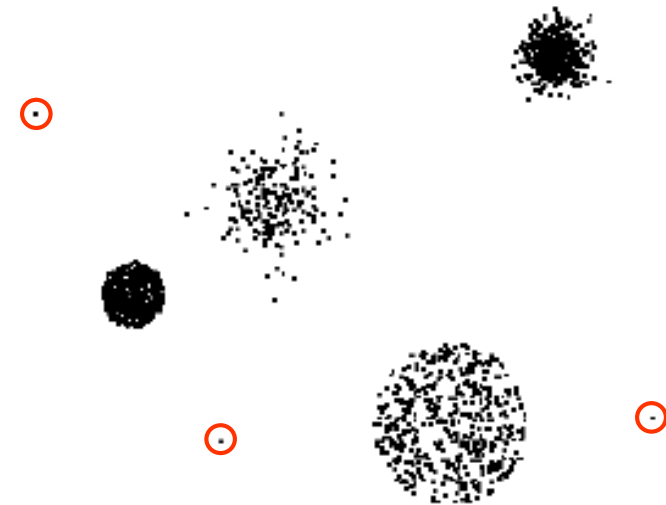
Association Rules



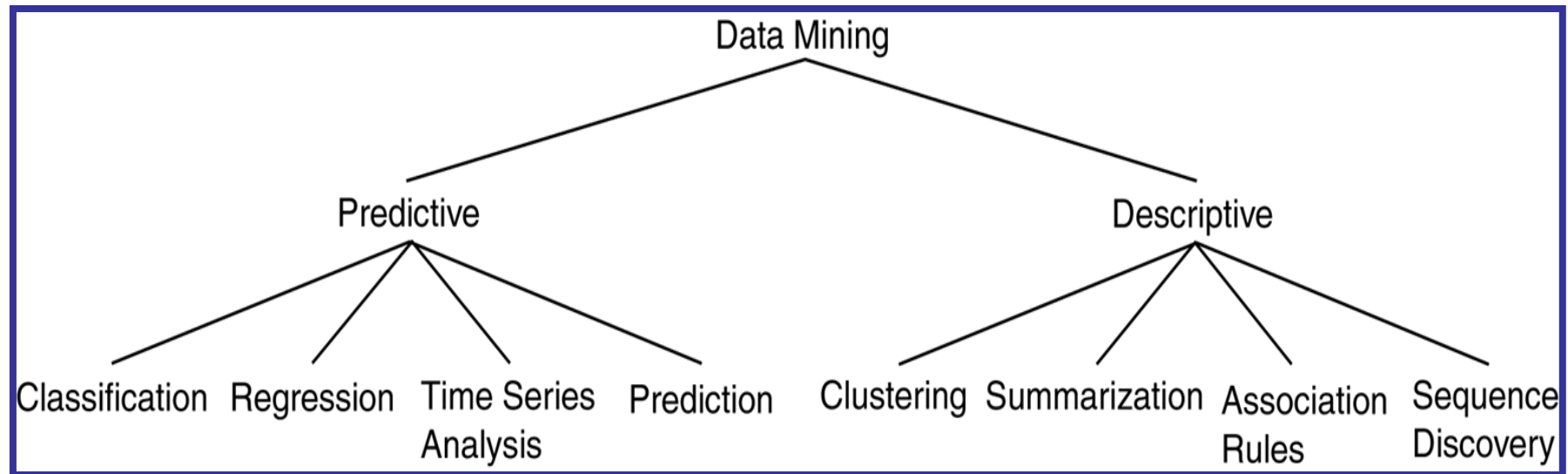
Predictive Modeling



Anomaly Detection

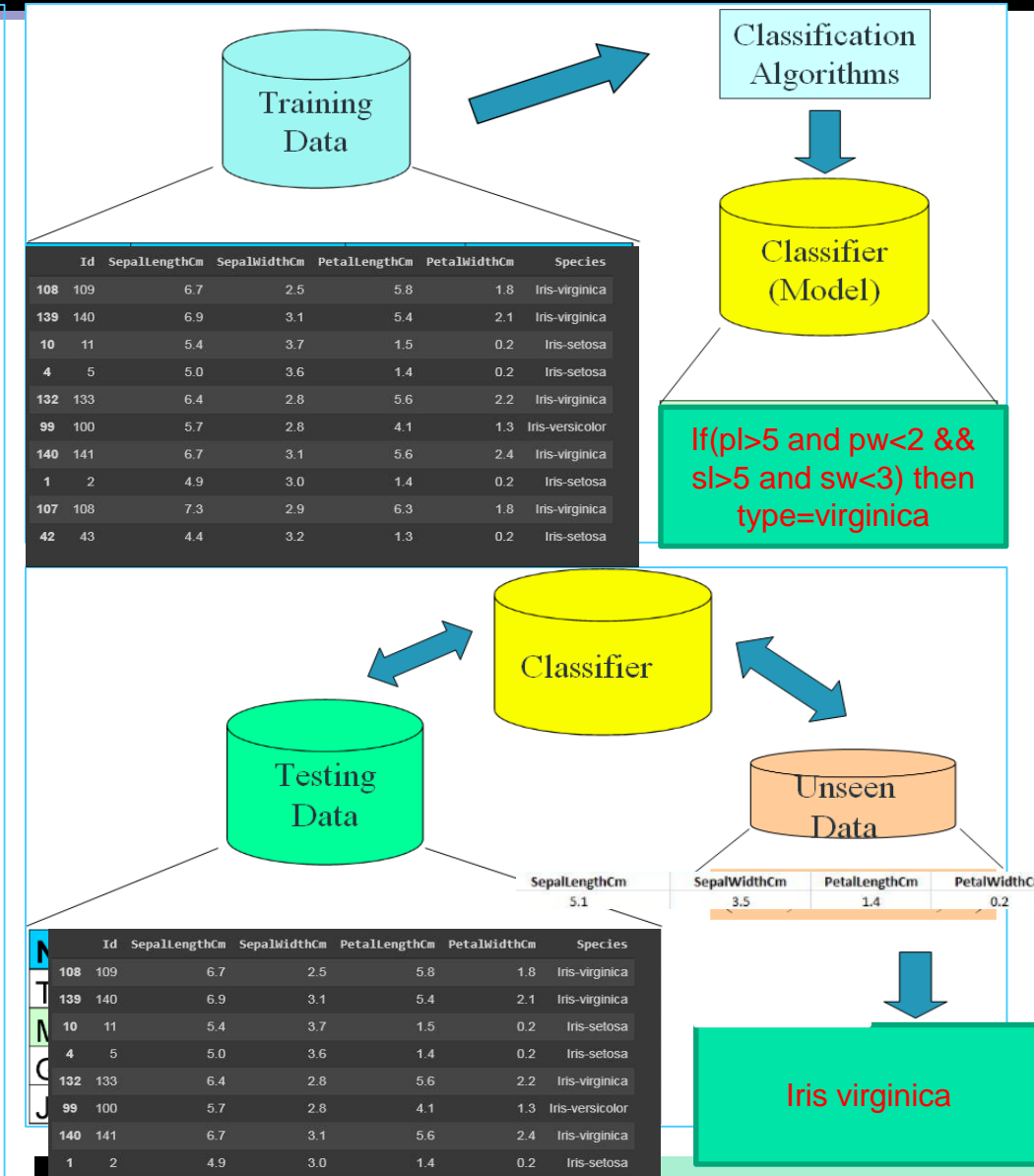


Data Mining Models and Tasks

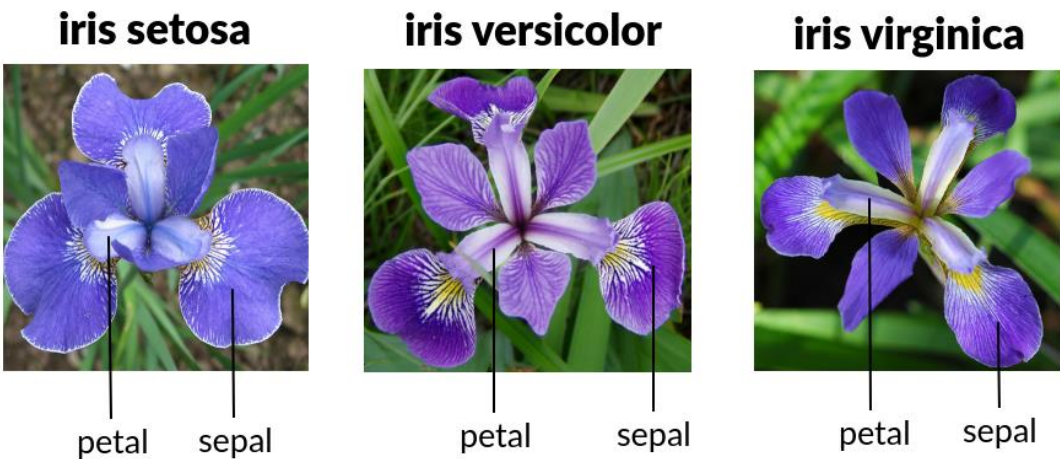


Supervised Learning

- Prediction methods are commonly referred to as supervised learning.
- Supervised methods are thought to attempt the discovery of the relationships between input attributes and a target attribute.
- A training set is given and the objective is to form a description that can be used to predict unseen examples.
- Methods:
 - **Classification**
 - The domain of the target attribute is finite and categorical.
 - A classifier must assign a class to a unseen example.
 - **Regression**
 - The target attribute is formed by infinite values.
 - To fit a model to learn the output target attribute as a function of input attributes.
 - **Time Series Analysis**
 - Making predictions in time.



Application-Feature-Class Label



	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
108	109	6.7	2.5	5.8	1.8	Iris-virginica
139	140	6.9	3.1	5.4	2.1	Iris-virginica
10	11	5.4	3.7	1.5	0.2	Iris-setosa
4	5	5.0	3.6	1.4	0.2	Iris-setosa
132	133	6.4	2.8	5.6	2.2	Iris-virginica
99	100	5.7	2.8	4.1	1.3	Iris-versicolor
140	141	6.7	3.1	5.6	2.4	Iris-virginica
1	2	4.9	3.0	1.4	0.2	Iris-setosa
107	108	7.3	2.9	6.3	1.8	Iris-virginica
42	43	4.4	3.2	1.3	0.2	Iris-setosa

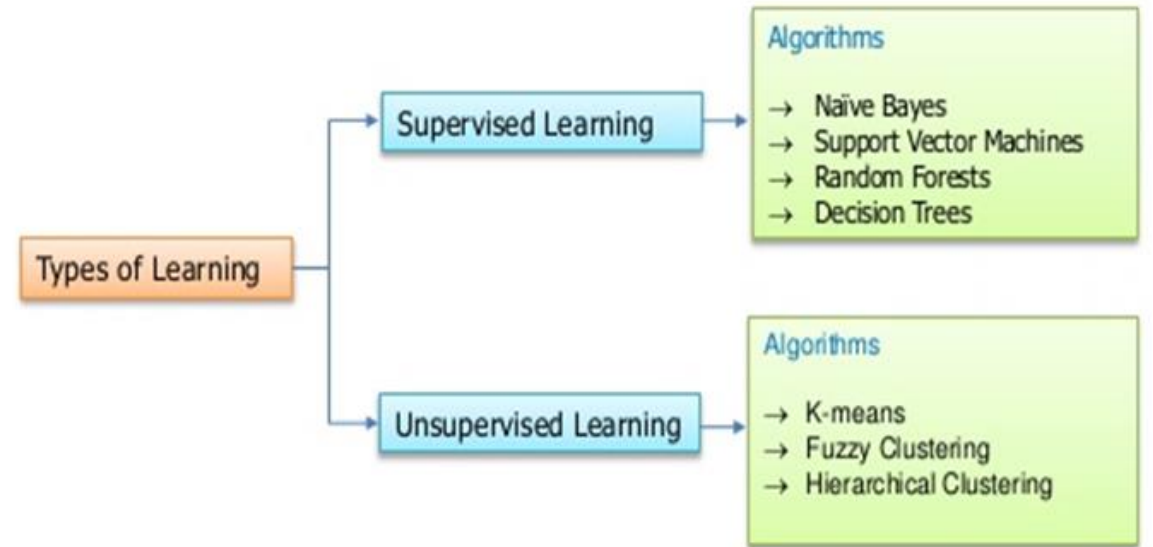
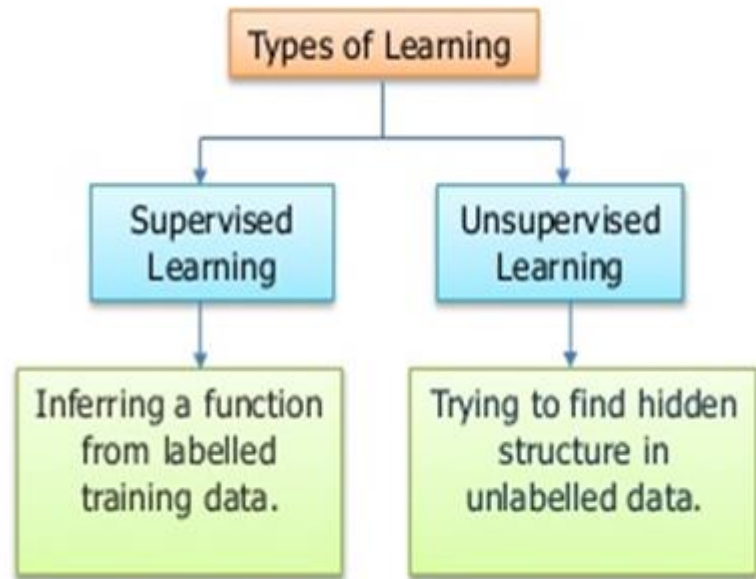
Unsupervised Learning

- There is no supervisor (labels) and only input data is available.
- The aim is to find regularities, irregularities, relationships, similarities and associations in the input.
- Methods:
 - Clustering
 - Association Rules
 - Pattern Mining
 - It is adopted as a more general term than frequent pattern mining or association mining.
 - Outlier Detection
 - Process of finding data examples with behaviours that are very different from the expectation (outliers or anomalies).

Semi Supervised Learning

- Semi-supervised learning is a class of machine learning tasks and techniques that also make use of unlabeled data for training – typically a small amount of labeled data with a large amount of unlabeled data
- Semi-supervised learning falls between unsupervised learning (without any labeled training data) and supervised learning (with completely labeled training data).

Unsupervised Learning Algorithms



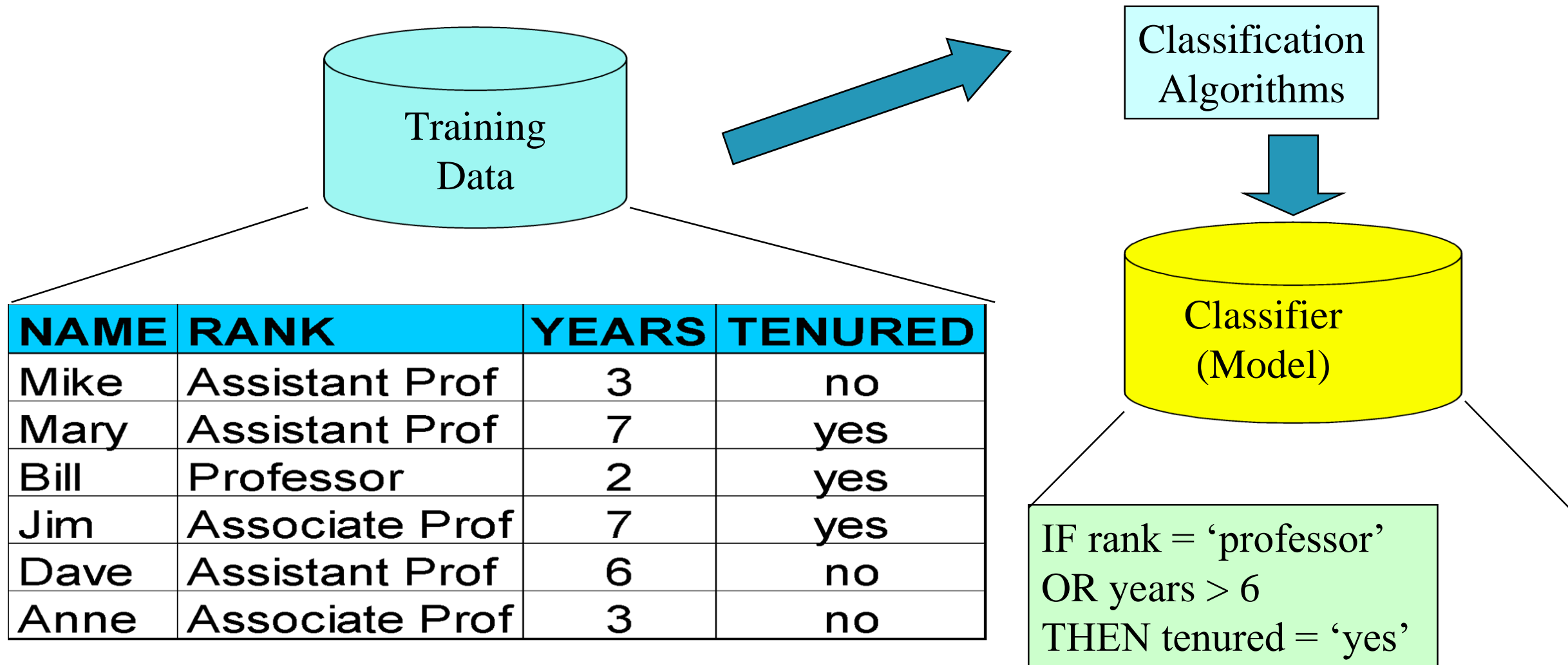
Classification

Classification maps data into predefined groups or classes

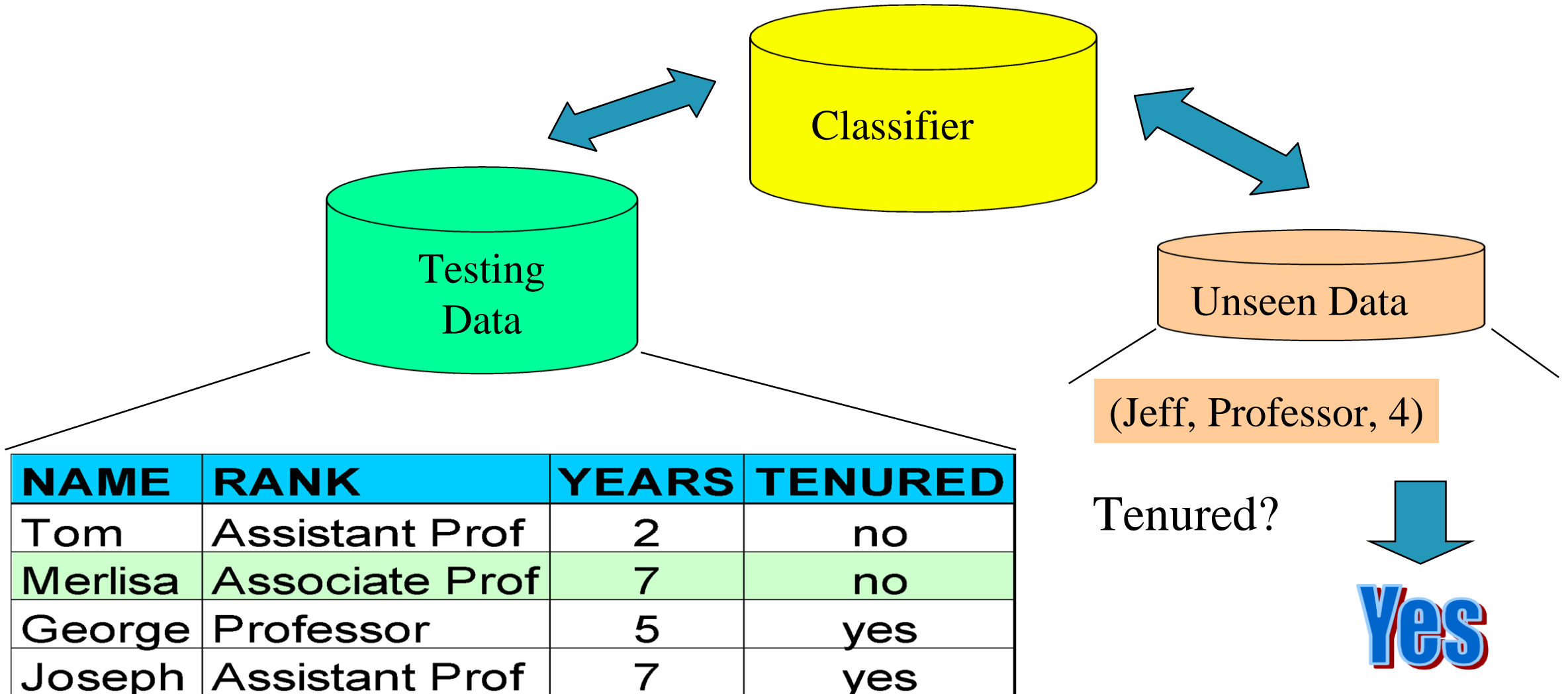
It is useful in

- Supervised learning
- Pattern recognition
- Prediction

Process (1): Model Construction



Process (2): Using the Model in Prediction



Classification Vs Prediction

- **Classification**

- Predicts categorical class labels (discrete or nominal)
- Classifies data (constructs a model) based on the training set and the values (**class labels**) in a classifying attribute and uses it in classifying new data

- **Prediction**

- Models continuous-valued functions, i.E., Predicts unknown or missing values

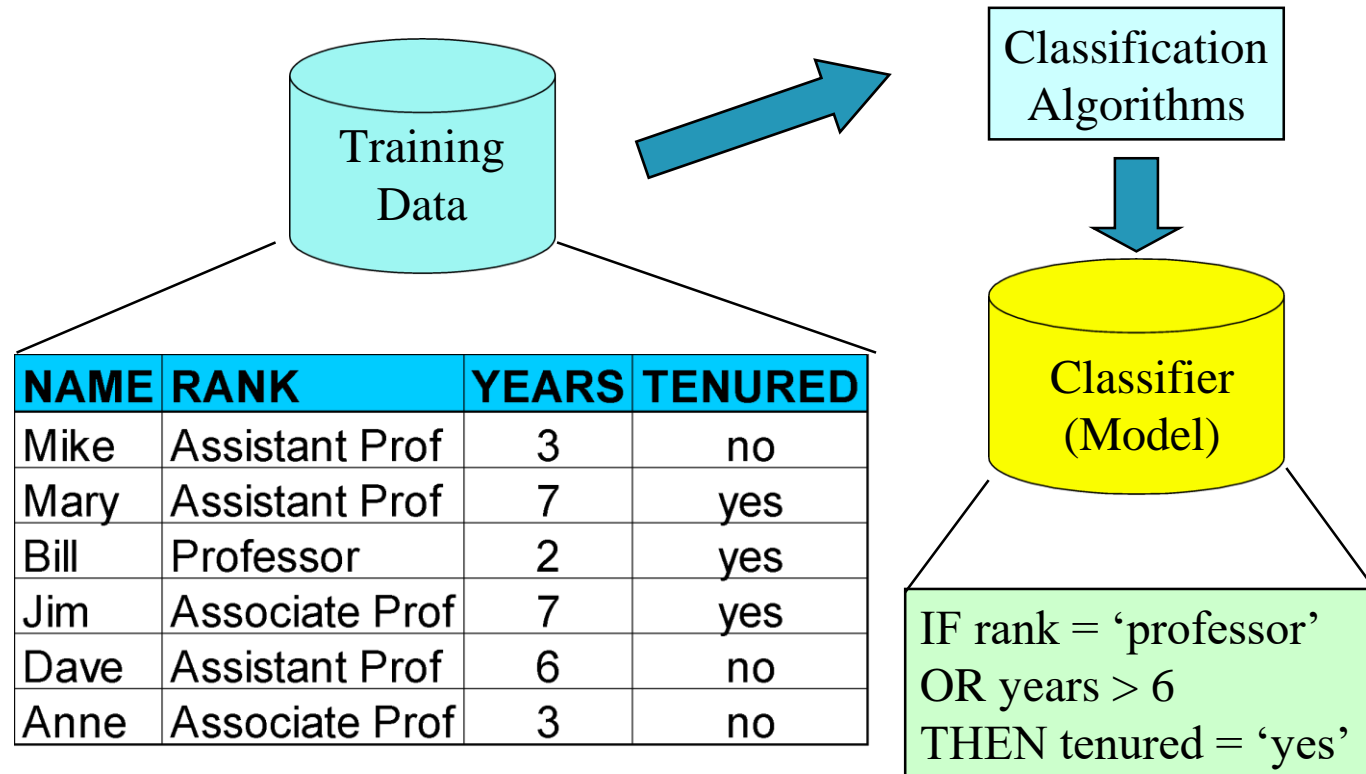
- **Typical applications**

- Credit approval : approve/not approve
- Target marketing : churn / not churn
- Medical diagnosis – disease : yes/no
- Fraud detection- legitimate/fraud

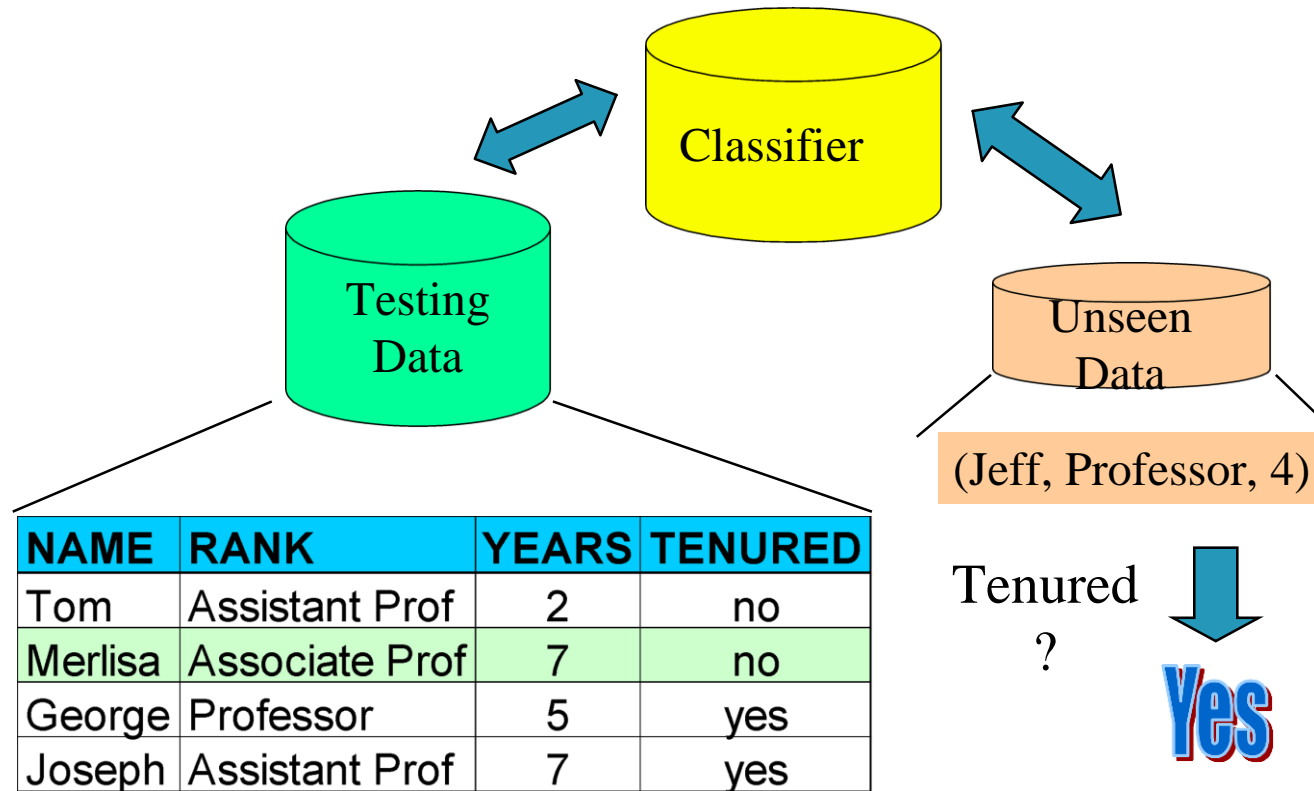
Classification Process

- **Model construction (Training):** describing a set of predetermined classes
 - Each tuple/sample is assumed to belong to a predefined class, as determined by the **class label attribute**
 - The set of tuples used for model construction is **training set**
 - The model is represented as classification rules, decision trees, or mathematical formulae
- **Model usage (Testing and Usage):** for classifying future or unknown objects
 - **Estimate accuracy** of the model
 - The known label of test sample is compared with the classified result from the model
 - Accuracy rate is the percentage of test set samples that are correctly classified by the model
 - Test set is independent of training set, otherwise over-fitting will occur
 - If the accuracy is acceptable, use the model to **classify data** tuples whose class labels are not known

Step 1: Model Construction



Step 2: Model Usage



Bayesian Classification

- A statistical classifier: performs *probabilistic prediction*, i.e., predicts class membership probabilities
- Foundation: Based on Bayes' Theorem.
- Performance: A simple Bayesian classifier, *naïve Bayesian classifier*, has comparable performance with decision tree and selected neural network classifiers
- Standard: Even when Bayesian methods are computationally intractable, they can provide a standard of optimal decision making against which other methods can be measured

Bayesian Classification - Example

age	income	student	credit rating	com
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Bayesian Classification - Probability

- *Probability*: How likely something is to happen

- Probability of an event happening =

Number of ways it can happen

Total number of outcomes

Bayesian Classification Theorem

- Let \mathbf{X} be a data sample (“*evidence*”): class label is unknown
- Let H be a *hypothesis* that X belongs to class C
- Classification is to determine $P(H|\mathbf{X})$, the probability that the hypothesis holds given the observed data sample \mathbf{X}
- $P(H)$ (*prior probability*), the initial probability
 - E.g., \mathbf{X} will buy computer, regardless of age, income, ...
- $P(\mathbf{X})$: probability that sample data is observed
- $P(\mathbf{X}|H)$ (*posteriori probability*), the probability of observing the sample \mathbf{X} , given that the hypothesis holds
 - E.g., Given that \mathbf{X} will buy computer, the prob. that X is 31..40, medium income

Bayesian Classification Theorem

- Given training data \mathbf{X} , *posteriori probability of a hypothesis* H , $P(H|\mathbf{X})$, follows the Bayes theorem

-

$$P(H|\mathbf{X}) = \frac{P(\mathbf{X}|H)P(H)}{P(\mathbf{X})}$$

- Informally, this can be written as
 - posteriori = likelihood x prior/evidence
- ***Predicts X belongs to C_i iff the probability $P(C_i|X)$ is the highest among all the $P(C_k|X)$ for all the k classes***
- Practical difficulty: require initial knowledge of many probabilities, significant computational cost

Naïve Bayesian Classification Theorem

- Let D be a training set of tuples and their associated class labels, and each tuple is represented by an n -D attribute vector $\mathbf{X} = (x_1, x_2, \dots, x_n)$
- Suppose there are m classes C_1, C_2, \dots, C_m .
- Classification is to derive the maximum posteriori, i.e., the maximal $P(C_i | \mathbf{X})$
- This can be derived from Bayes' theorem

$$P(C_i | \mathbf{X}) = \frac{P(\mathbf{X} | C_i)P(C_i)}{P(\mathbf{X})}$$

- Since $P(\mathbf{X})$ is constant for all classes, only

$$P(C_i | \mathbf{X}) = P(\mathbf{X} | C_i)P(C_i)$$

- needs to be maximized

Naïve Bayesian Classification - Example

Class:

C1:buys_computer = 'yes'

C2:buys_computer = 'no'

Data sample

X = (age <=30,

Income = medium,

Student = yes

Credit_rating = Fair)

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Naïve Bayesian Classification - Example

Test for $X = (\text{age} \leq 30, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit_rating} = \text{fair})$

- $P(C_i)$: $P(\text{buys_computer} = \text{"yes"}) = 9/14 = 0.643$ $P(\text{buys_computer} = \text{"no"}) = 5/14 = 0.357$

Compute $P(X|C_i)$ for each class

$$P(\text{age} = \text{"<=30"} \mid \text{buys_computer} = \text{"yes"}) = 2/9 = 0.222$$

$$P(\text{age} = \text{"<= 30"} \mid \text{buys_computer} = \text{"no"}) = 3/5 = 0.6$$

$$P(\text{income} = \text{"medium"} \mid \text{buys_computer} = \text{"yes"}) = 4/9 = 0.444$$

$$P(\text{income} = \text{"medium"} \mid \text{buys_computer} = \text{"no"}) = 2/5 = 0.4$$

$$P(\text{student} = \text{"yes"} \mid \text{buys_computer} = \text{"yes"}) = 6/9 = 0.667$$

$$P(\text{student} = \text{"yes"} \mid \text{buys_computer} = \text{"no"}) = 1/5 = 0.2$$

$$P(\text{credit_rating} = \text{"fair"} \mid \text{buys_computer} = \text{"yes"}) = 6/9 = 0.667$$

$$P(\text{credit_rating} = \text{"fair"} \mid \text{buys_computer} = \text{"no"}) = 2/5 = 0.4$$

- $P(X|C_i)$: $P(X \mid \text{buys_computer} = \text{"yes"}) = 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$

$$P(X \mid \text{buys_computer} = \text{"no"}) = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$$

- $P(X|C_i) \cdot P(C_i)$: $P(X \mid \text{buys_computer} = \text{"yes"}) \cdot P(\text{buys_computer} = \text{"yes"}) = \mathbf{0.028}$

$$P(X \mid \text{buys_computer} = \text{"no"}) \cdot P(\text{buys_computer} = \text{"no"}) = \mathbf{0.007}$$

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

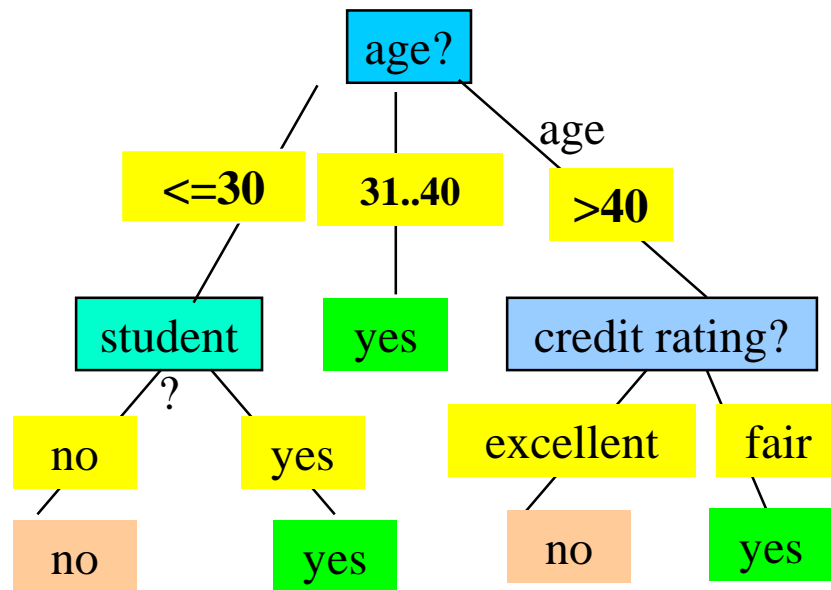
0.028 > 0.007 ..

**Therefore, X belongs to class
("buys_computer = yes")**

Comment on Naive Bayes Classification

- Advantages
 - Easy to implement
 - Good results obtained in most of the cases
- Disadvantages
 - Assumption: class conditional independence, therefore loss of accuracy

Decision Tree Classification



age	income	student	credit rating	buys computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Decision Tree Terminologies

- Decision tree may be n -ary, $n \geq 2$.
- There is a special node called **root node**.
- **Internal nodes** are test attribute/decision attribute
- **leaf nodes** are class labels
- Edges of a node represent the **outcome for a value** of the test node.
- In a path, a node with same label **is never repeated**.
- Decision tree **is not unique**, as different ordering of internal nodes can give different decision tree

Rule Extraction From decision Tree

- Rules are easier to understand than large trees
- One rule is created for each path from the root to a leaf
- Each attribute-value pair along a path forms a conjunction: the leaf holds the class prediction
- Rules are mutually exclusive and exhaustive

Example: Rule extraction from our *buys_computer* decision-tree

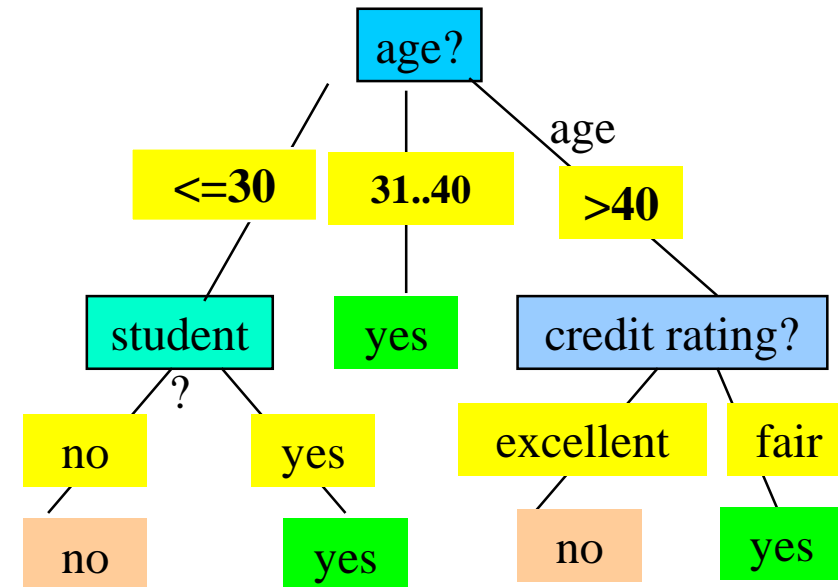
IF *age* = young AND *student* = no THEN *buys_computer* = no

IF *age* = young AND *student* = yes THEN *buys_computer* = yes

IF *age* = mid-age THEN *buys_computer* = yes

IF *age* = old AND *credit_rating* = excellent THEN *buys_computer* = yes

IF *age* = young AND *credit_rating* = fair THEN *buys_computer* = no



Decision Tree Algorithm

- **Basic algorithm (a greedy algorithm)**
 - Tree is constructed in a **top-down recursive divide-and-conquer manner**
 - At start, all the training examples are at the root
 - Attributes are categorical (if continuous-valued, they are discretized in advance)
 - Examples are partitioned recursively based on selected attributes
 - Test attributes are selected on the basis of a heuristic or statistical measure (e.g., **information gain**)
- **Conditions for stopping partitioning**
 - All samples for a given node belong to the same class
 - There are no remaining attributes for further partitioning – **majority voting** is employed for classifying the leaf
 - There are no samples left

Decision Tree Algorithm

Test Attribute Selection

- **Select the attribute with the highest information gain**

- Let p_i be the probability that an arbitrary tuple in D belongs to class C_i , estimated by $|C_{i,D}|/|D|$

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

- **Expected information** (entropy) needed to classify a tuple in D: $Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times I(D_j)$

- **Information** needed (after using A to split D into v partitions) to classify D:

$$Gain(A) = Info(D) - Info_A(D)$$

- **Information gained** by branching on attribute A

Attribute Selection : Information gain

Class P: buys_computer = "yes"

Class N: buys_computer = "no"

1 : Calculate Entropy for Class Labels

$$Info(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940$$

2 : Calculate Information of Each Attribute (one by one)

age	p _i	n _i	I(p _i , n _i)
<=30	2	3	0.971
31...40	4	0	0
>40	3	2	0.971

$$Info_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$$

$\frac{5}{14} I(2,3)$ means "age <=30" has 5 out of 14 samples, with 2 yes'es and 3 no's.

age	income	student	credit rating	buys computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

3 : Calculate Gain of Each Attribute (one by one)

$$Gain(age) = Info(D) - Info_{age}(D) = 0.246$$

$$Gain(income) = 0.029$$

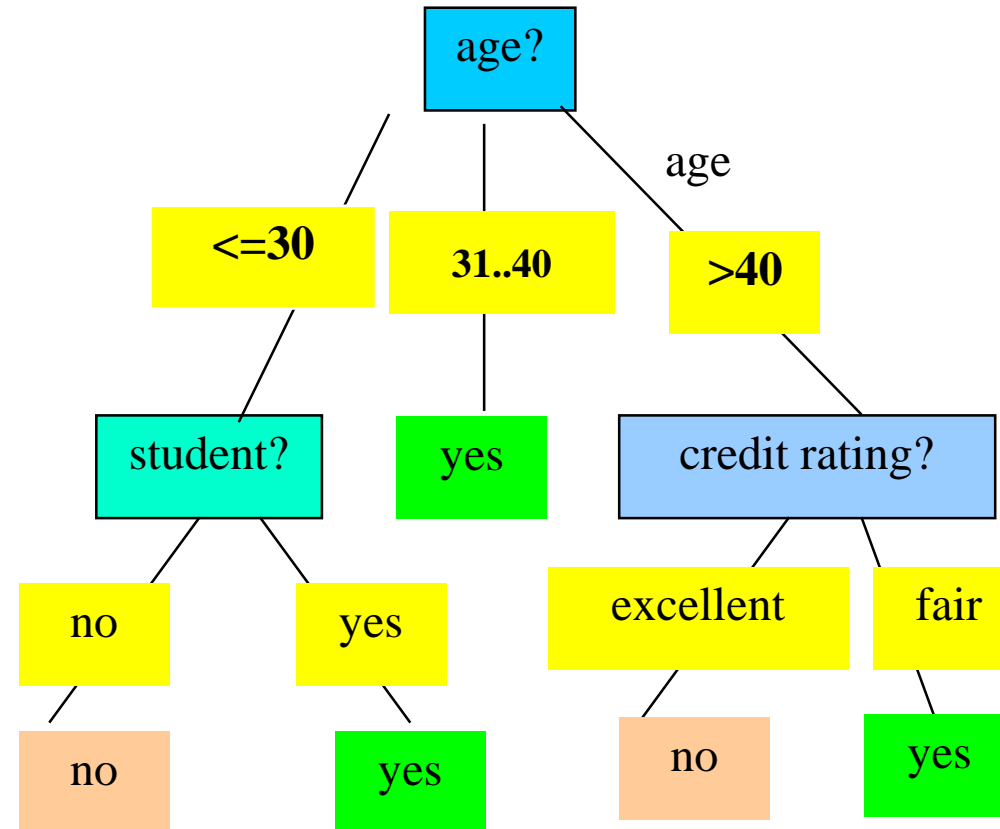
$$Gain(student) = 0.151$$

$$Gain(credit_rating) = 0.048$$

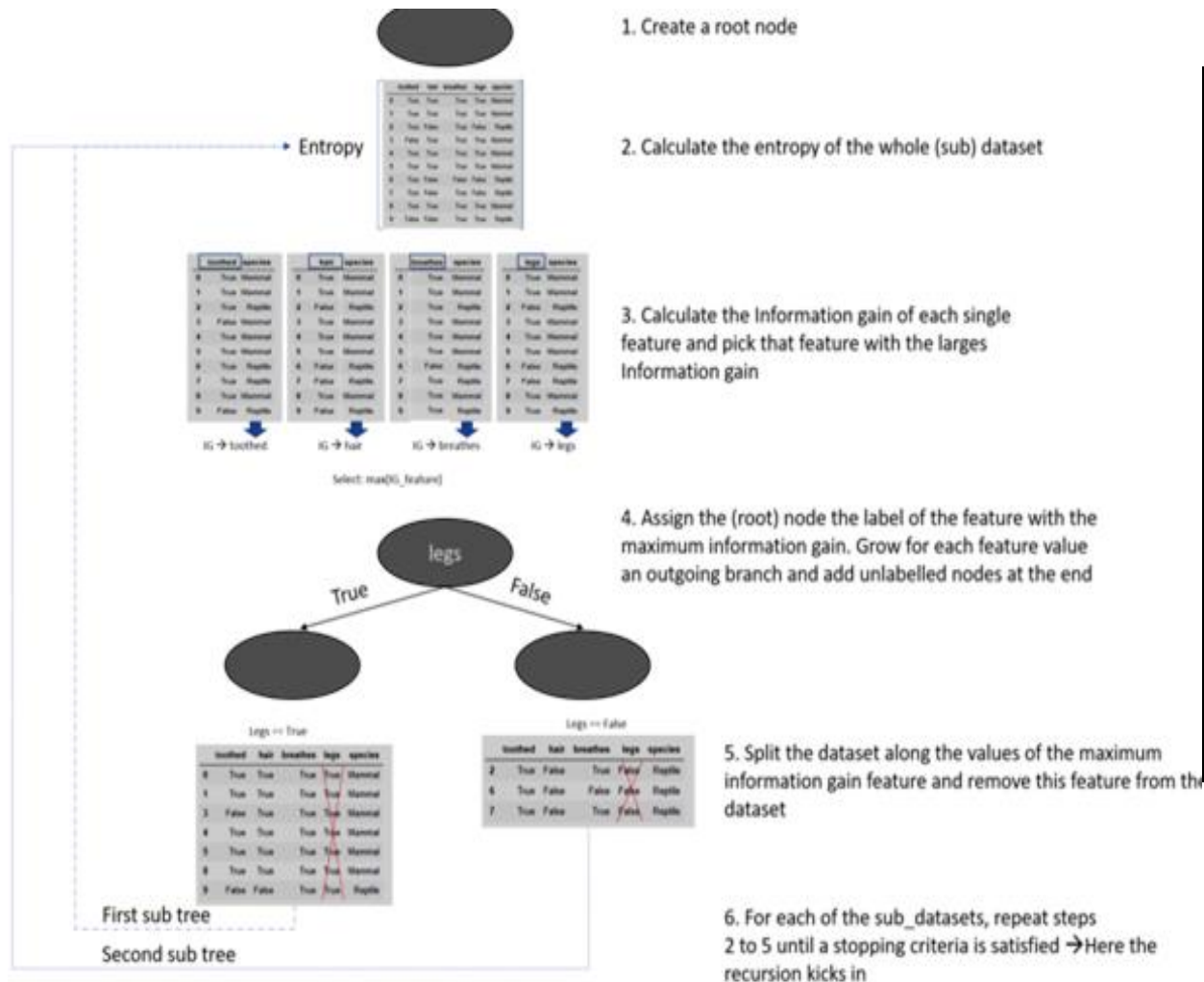
4 : Select Attribute with max gain as a root attribute

Gain(age) > any other gain. so age is selected as root node

Final DT of the buys_computer dataset

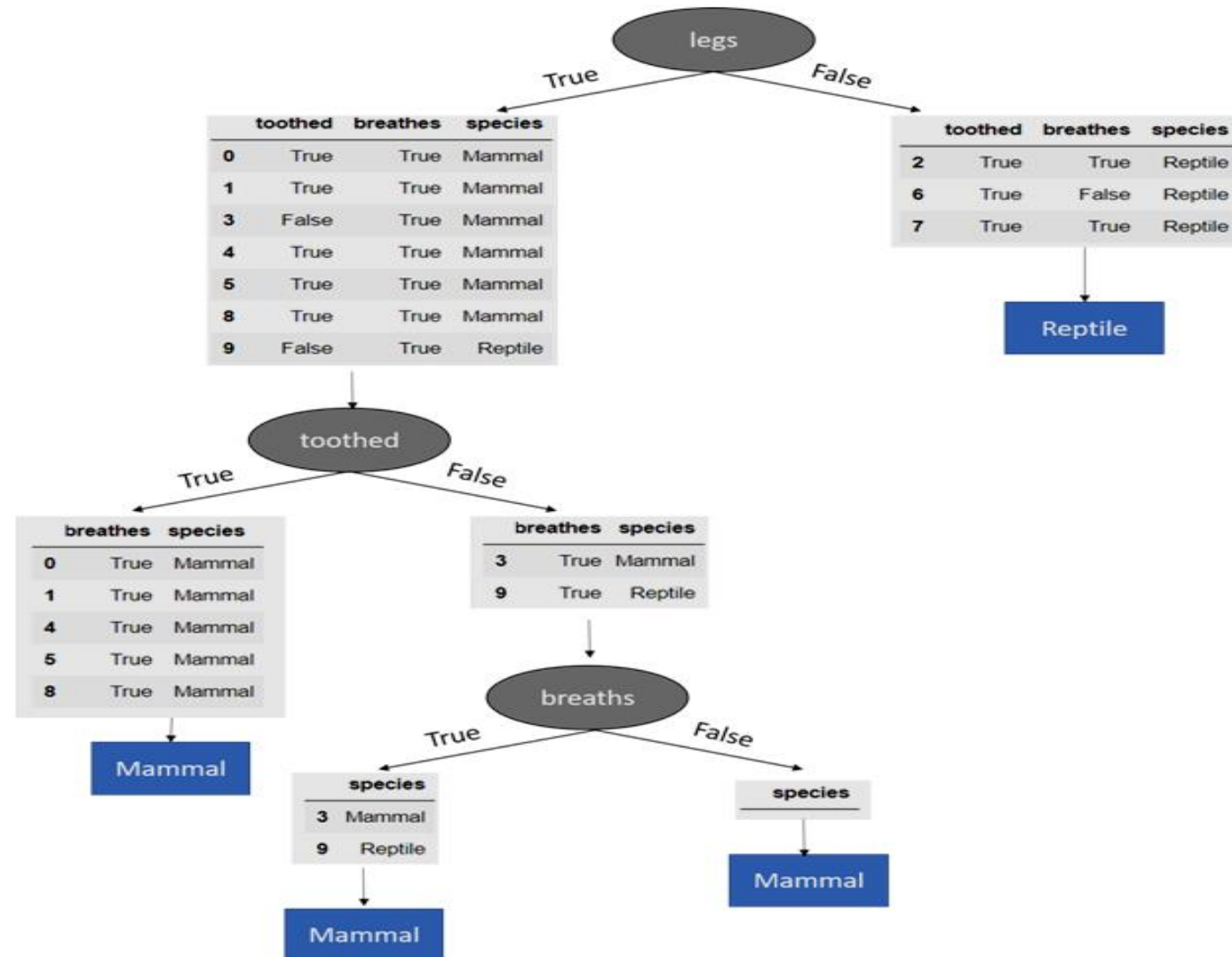


Example : Usage of Information Gain and Entropy in DT Creation



	toothed	hair	breathes	legs	species
0	True	True	True	True	Mammal
1	True	True	True	True	Mammal
2	True	False	True	False	Reptile
3	False	True	True	True	Mammal
4	True	True	True	True	Mammal
5	True	True	True	True	Mammal
6	True	False	False	False	Reptile
7	True	False	True	False	Reptile
8	True	True	True	True	Mammal
9	False	False	True	True	Reptile

Decision Tree Algorithm- Example



Classifier Evaluation Metrics: Confusion Matrix

Confusion Matrix:

Actual class \ Predicted class	C_1	$\neg C_1$
C_1	True Positives (TP)	False Negatives (FN)
$\neg C_1$	False Positives (FP)	True Negatives (TN)

Example of Confusion Matrix:

Actual class \ Predicted class	buy_computer = yes	buy_computer = no	Total
buy_computer = yes	6954	46	7000
buy_computer = no	412	2588	3000
Total	7366	2634	10000

- Given m classes, an entry, $\mathbf{CM}_{i,j}$ in a **confusion matrix** indicates # of tuples in class i that were labeled by the classifier as class j
- May have extra rows/columns to provide totals

Classifier Evaluation Metrics: Accuracy, Error Rate, Sensitivity and Specificity

A\P	C	¬C	
C	TP	FN	P
¬C	FP	TN	N
	P'	N'	All

- **Classifier Accuracy**, or recognition rate: percentage of test set tuples that are correctly classified

$$\text{Accuracy} = (TP + TN) / \text{All}$$

- **Error rate**: $1 - \text{accuracy}$, or
Error rate = $(FP + FN) / \text{All}$

- **Class Imbalance Problem:**

- One class may be *rare*, e.g. fraud, or HIV-positive
- Significant *majority of the negative class* and minority of the positive class

- **Sensitivity**: True Positive recognition rate

- **Sensitivity** = TP / P

- **Specificity**: True Negative recognition rate

- **Specificity** = TN / N

Classifier Evaluation Metrics: Precision and Recall, and F-measures

- **Precision:** exactness – what % of tuples that the classifier labeled as positive are actually positive

$$\text{precision} = \frac{TP}{TP + FP}$$

- **Recall:** completeness – what % of positive tuples did the classifier label as positive?

$$\text{recall} = \frac{TP}{TP + FN}$$

- Perfect score is 1.0
- Inverse relationship between precision & recall
- **F measure (F_1 or F-score):** harmonic mean of precision and recall,

$$F = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

- **F_β :** weighted measure of precision and recall
 - assigns β times as much weight to recall as to precision

$$F_\beta = \frac{(1 + \beta^2) \times \text{precision} \times \text{recall}}{\beta^2 \times \text{precision} + \text{recall}}$$

Classifier Evaluation Metrics: Example

Actual Class\Predicted class	cancer = yes	cancer = no	Total	Recognition(%)
cancer = yes	90	210	300	30.00 (<i>sensitivity</i>)
cancer = no	140	9560	9700	98.56 (<i>specificity</i>)
Total	230	9770	10000	96.40 (<i>accuracy</i>)

■ $Precision = 90/230 = 39.13\%$

$Recall = 90/300 = 30.00\%$

Summary- Classification Algorithm

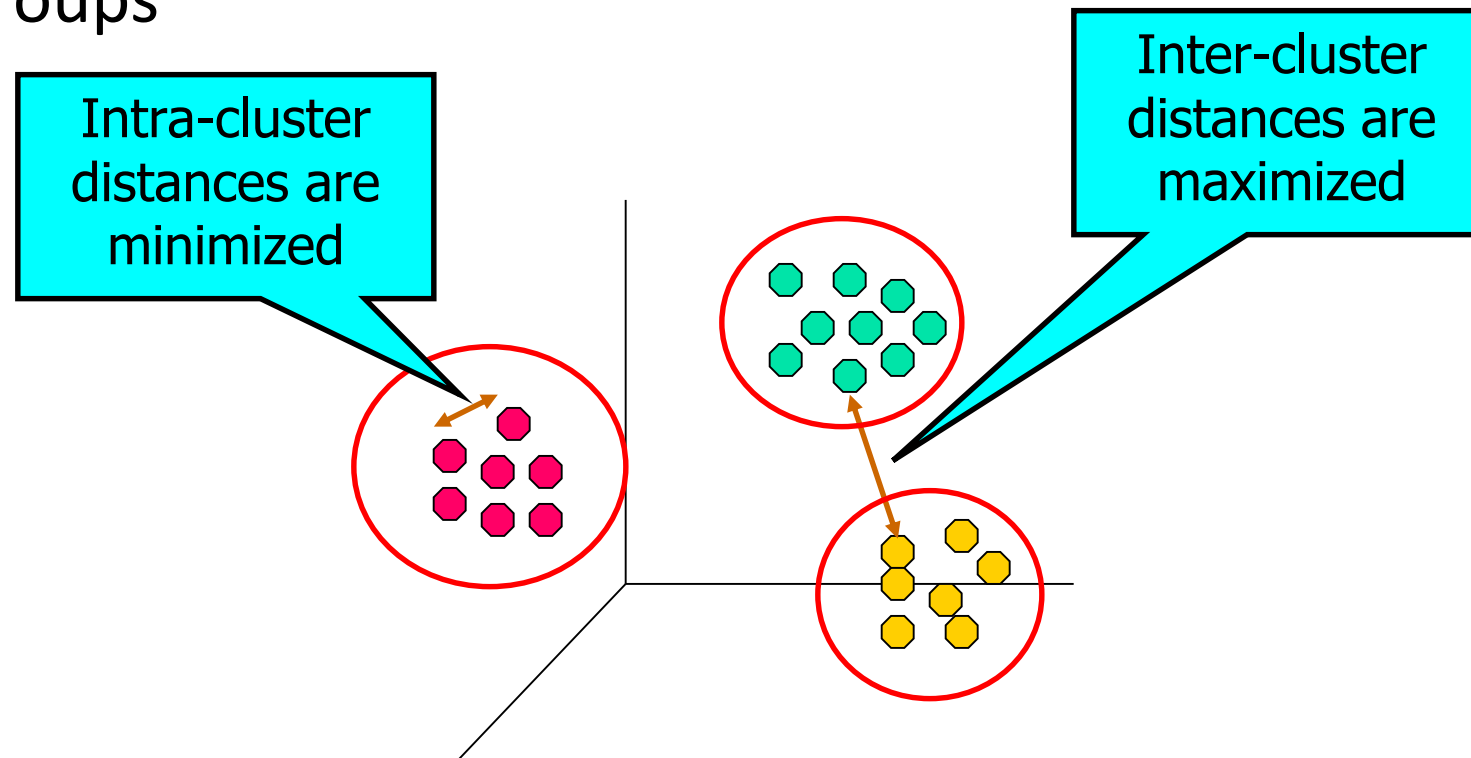
- **Classification** is a form of data analysis that extracts **models** describing important data classes.
- Effective and scalable methods have been developed for **decision tree induction, Naive Bayesian classification, rule-based classification**, and many other classification methods.

Clustering

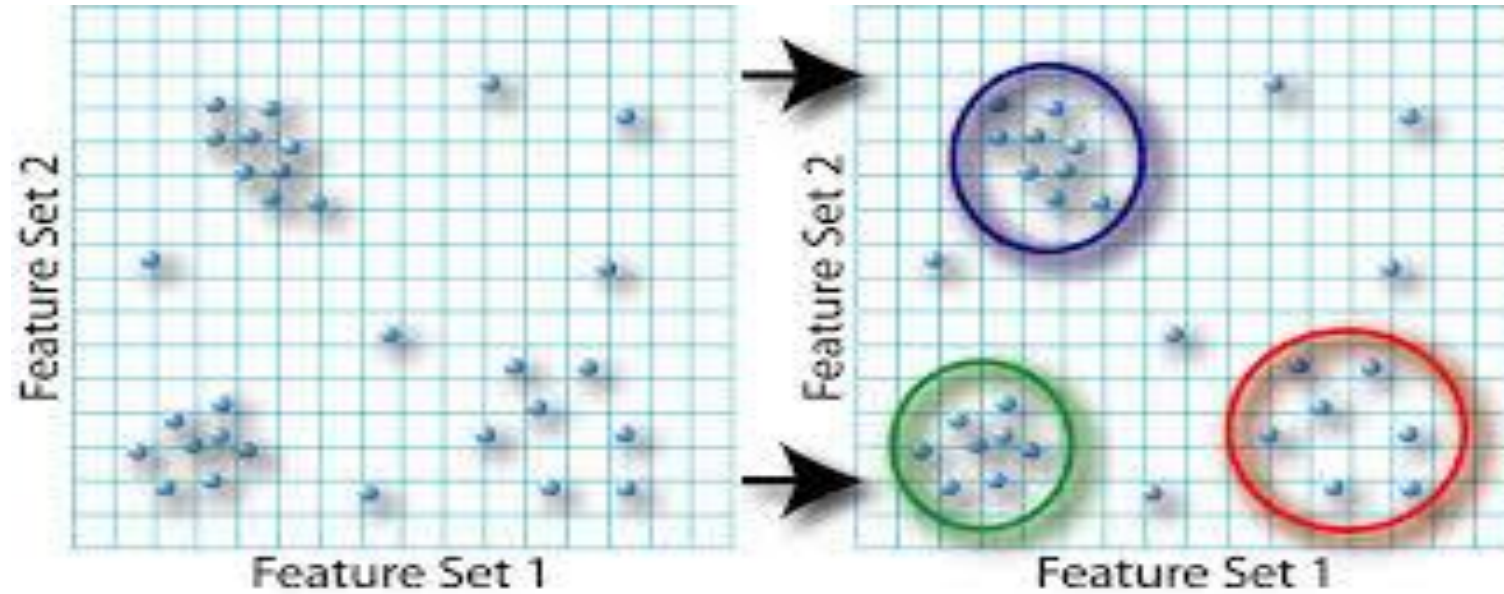
- Organizing data into classes such that there is
 - high intra-class similarity
 - low inter-class similarity
- Finding the class labels and the number of classes directly from the data (in contrast to classification).
- More informally, finding natural groupings among objects.
- Also called unsupervised learning, sometimes called classification by statisticians and sorting by psychologists and segmentation by people in marketing

Definition

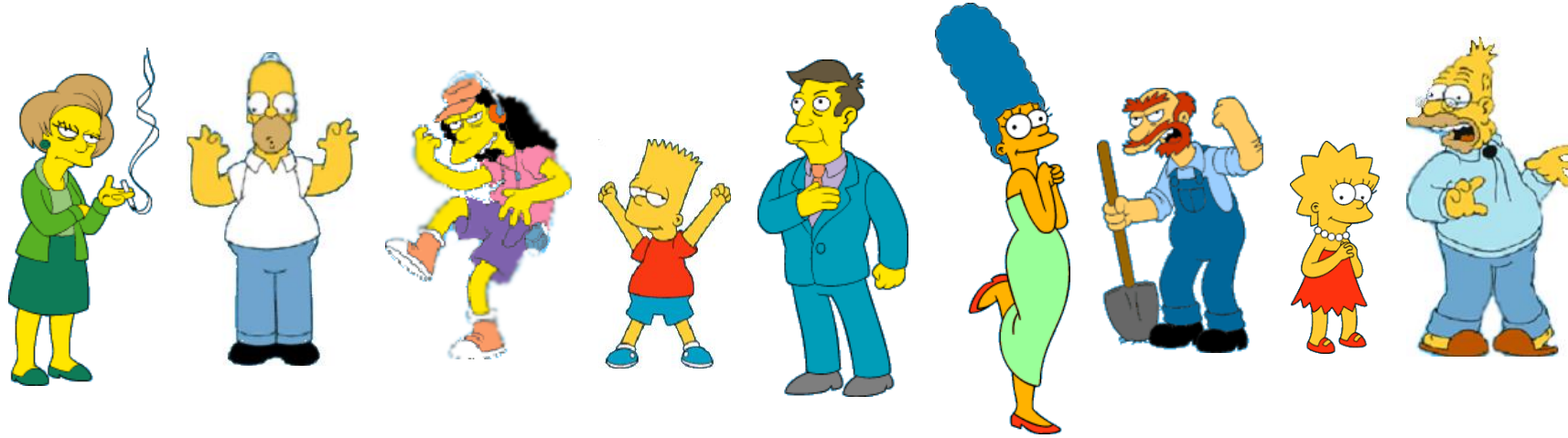
- Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups



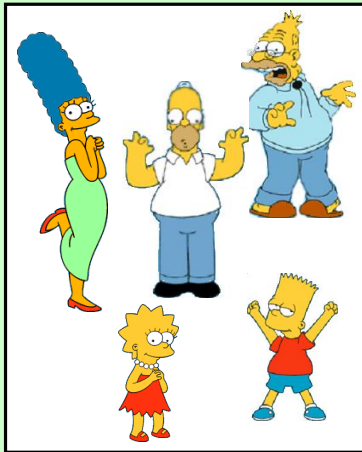
Clustering Example



What is a natural grouping among these objects?



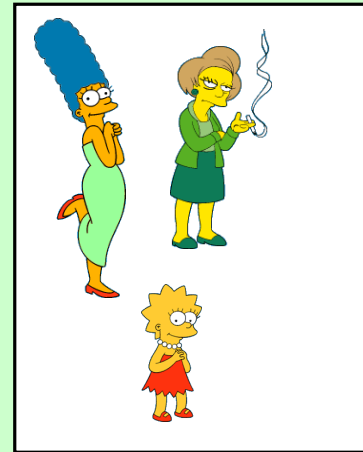
Clustering is subjective



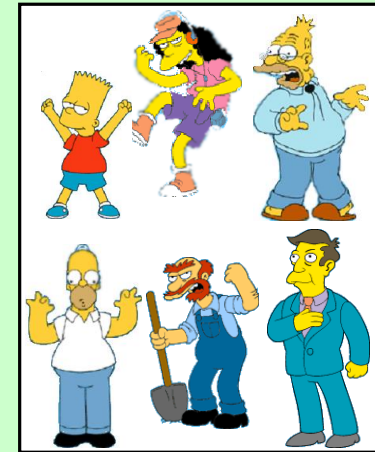
Simpson's Family



School Employees



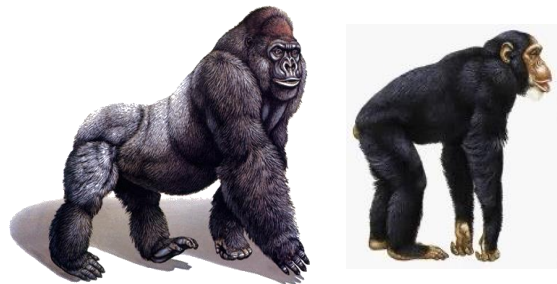
Females



Males

Defining Distance Measures

Definition: Let O_1 and O_2 be two objects from the universe of possible objects. The distance (dissimilarity) between O_1 and O_2 is a real number denoted by $D(O_1, O_2)$



0.23

Peter

Piotr



3



342.7

Common Distance measures:

Distance measure will determine how the *similarity* of two elements is calculated and it will influence the shape of the clusters.

They include:

1. The Euclidean distance (also called 2-norm distance) is given by:

$$d(x, y) = \sqrt{\sum_{i=1}^P |x_i - y_i|^2}$$

1. The Manhattan distance (also called taxicab norm or 1-norm) is given by:

$$d(x, y) = \sum_{i=1}^P |x_i - y_i|$$

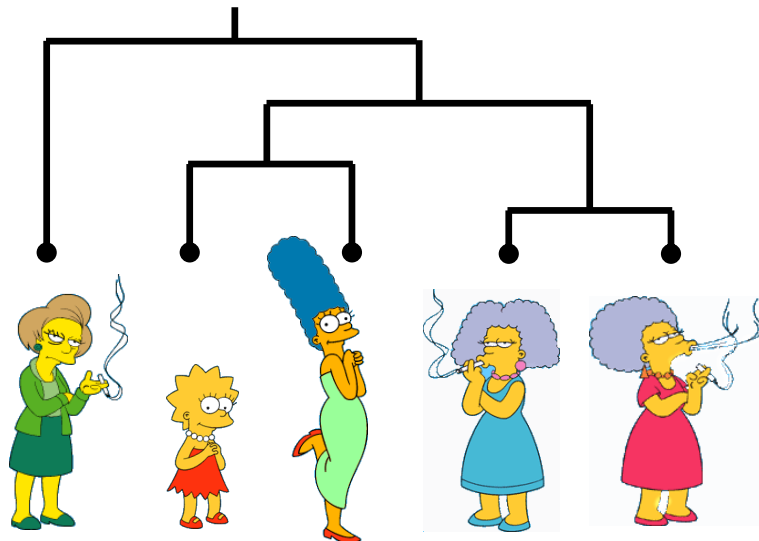
Types of Clustering

- Hierarchical clustering(BIRCH)
 - A set of nested clusters organized as a hierarchical tree
- Partitional Clustering(k-means,k-mediods)
 - A division data objects into non-overlapping (distinct) subsets (i.e., clusters) such that each data object is in exactly one subset
- Density – Based(DBSCAN)
 - Based on density functions
- Grid-Based(STING)
 - Based on multiple-level granularity structure
- Model-Based(SOM)
 - Hypothesize a model for each of the clusters and find the best fit of the data to the given model

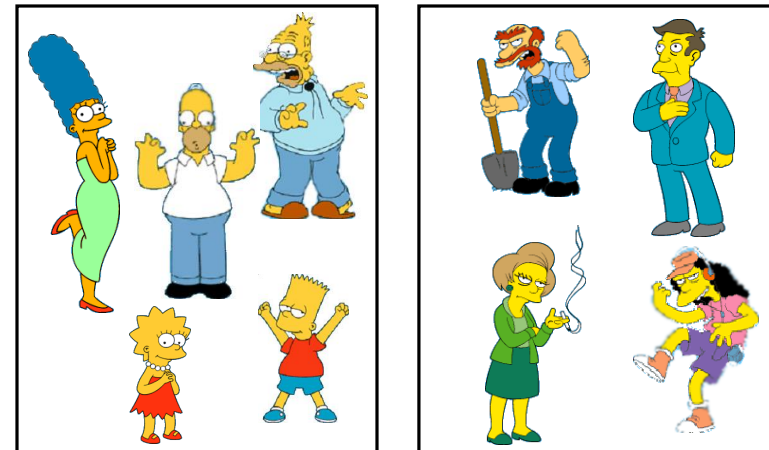
Two types of Clustering

- **Partitional algorithms:** Construct various partitions and then evaluate them by some criterion (we will see an example called BIRCH)
- **Hierarchical algorithms:** Create a hierarchical decomposition of the set of objects using some criterion

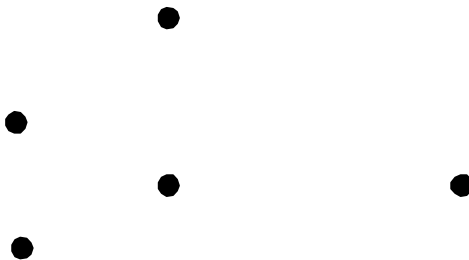
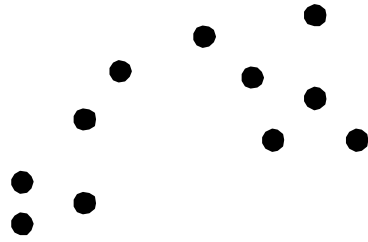
Hierarchical



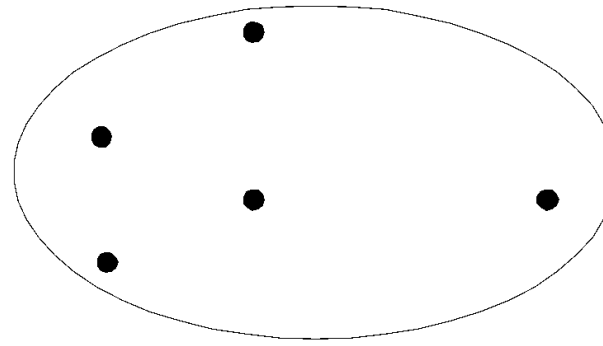
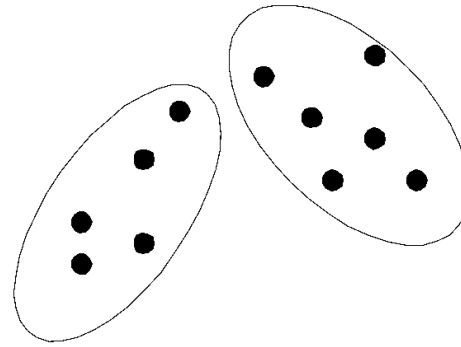
Partitional



Partitional Clustering



Original Points



A Partitional Clustering

Desirable Properties of Clustering

- Scalability (in terms of both time and space)
- Ability to deal with different data types
- Minimal requirements for domain knowledge to determine input parameters
- Able to deal with noise and outliers
- Insensitive to order of input records
- Incorporation of user-specified constraints
- Interpretability and usability

K-MEANS CLUSTERING

- The **k-means algorithm** is an algorithm to cluster n objects based on attributes into k partitions, where $k < n$.
- It assumes that the object attributes form a vector space.
- An algorithm for partitioning (or clustering) N data points into K disjoint subsets S_j containing data points so as to minimize the sum-of-squares criterion

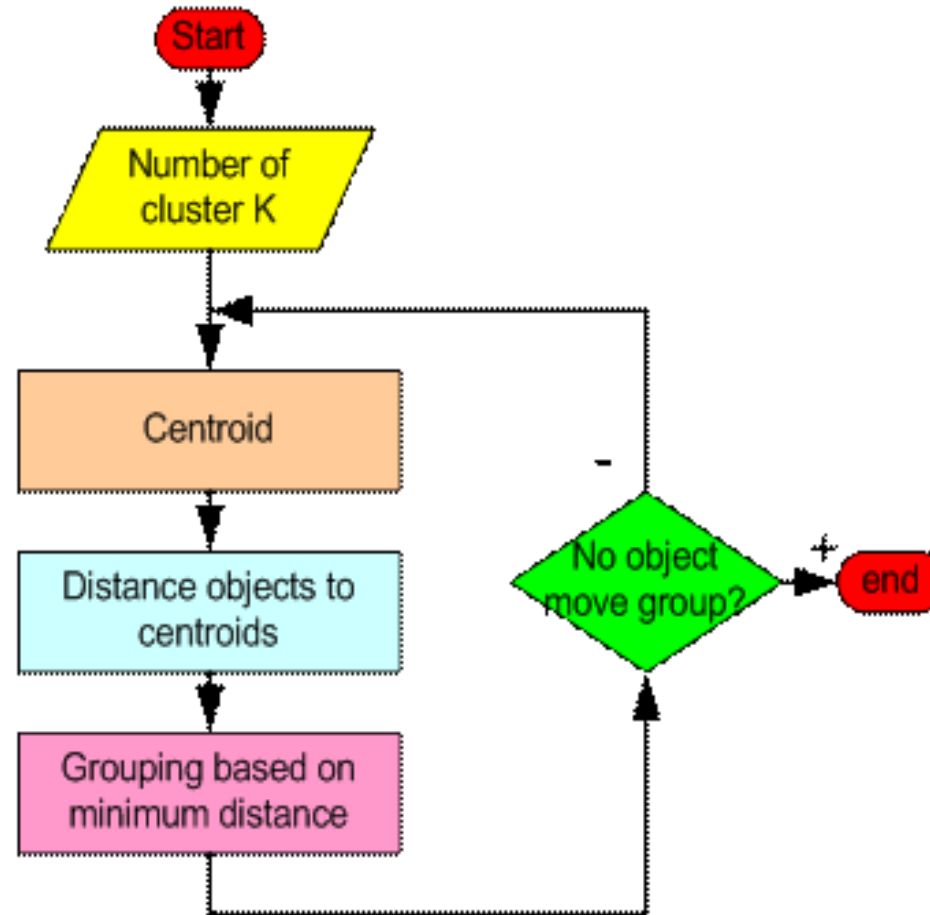
$$J = \sum_{j=1}^K \sum_{n \in S_j} |x_n - \mu_j|^2,$$

- Where x_n is a vector representing the n^{th} data point and μ_j is the geometric centroid of the data points in S_j .

K-MEANS CLUSTERING

- Simply speaking k-means clustering is an algorithm to classify or to group the objects based on attributes/features into **K number of group**.
- K is positive integer number.
- The grouping is done by minimizing the **sum of squares of distances** between data and the corresponding **cluster centroid**.

Working of K-MEANS CLUSTERING



Classical Partitioning Method- K mean

- First, it randomly selects K of the objects, each of which initially means represents cluster mean or center.
- For each of the remaining objects, an object is assigned to the cluster to which it is most similar, based on the distance between the object and the cluster mean.
- It then computes the new mean for each cluster.
- This process iterates until the criterion function converges.
- Typically, the Square error criterion is used, defines as

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

- **E**- Sum of square error for all objects in the data set
- **P** is the point in space representing given object;
- **mi**- is the mean of cluster Ci

Classical Partitioning Method- K mean

- In other words, for each object in each cluster , the distance from the object to its cluster center is squared, and the distances are summed.
- This criterion tries to make the resulting k clusters as compact and as separate as possible.

Working of K-MEANS CLUSTERING

- Begin with a decision on the value of **k = number of clusters**
- Arbitrarily assign k objects from D as the initial cluster centers
- Each object is distributed to a cluster based on the cluster center to which it is the nearest.
- Next, the cluster centers are updated i.e. mean value of each cluster is recalculated based on the current objects in the cluster
- Using the new cluster centers, the objects are redistributed to the clusters based on which cluster center is the nearest.
- This process iterates
- Eventually, no redistribution of the objects in any occurs, and so the process terminates
- Resulting clusters are returned by the clustering process.

Algorithm of K-MEANS CLUSTERING

Algorithm: k -means. The k -means algorithm for partitioning, where each cluster's center is represented by the mean value of the objects in the cluster.

Input:

- k : the number of clusters,
- D : a data set containing n objects.

Output: A set of k clusters.

Method:

- (1) arbitrarily choose k objects from D as the initial cluster centers;
- (2) **repeat**
- (3) (re)assign each object to the cluster to which the object is the most similar,
 based on the mean value of the objects in the cluster;
- (4) update the cluster means, that is, calculate the mean value of the objects for
 each cluster;
- (5) **until** no change;

Data = {1,2,3,4,5,6,7,8,9} K =2 with initial centroid C1=3 C2=7



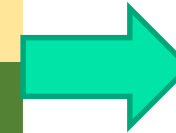
Calculate Manhattan Distance of each Data Object with Centroids

	C1=3	C2=7
1	2	6
2	1	5
3	0	4
4	1	3
5	2	2
6	3	1
7	4	0
8	5	1
9	6	2



newC1=1+2+3+4/4
new C2=5+6+7+8+9/5

	C1=3	C2=7	
1	2	6	C1
2	1	5	C1
3	0	4	C1
4	1	3	C1
5	2	2	C2
6	3	1	C2
7	4	0	C2
8	5	1	C2
9	6	2	C2



newC1=1+2+3+4/4
new C2=5+6+7+8+9/5

2.5
7

	newC1=2.5	C2=7	
1	1.5	6	C1
2	0.5	5	C1
3	0.5	4	C1
4	1.5	3	C1
5	2.5	2	C2
6	3.5	1	C2
7	4.5	0	C2
8	5.5	1	C2
9	6.5	2	C2

STOP!! As same cluster obtained in next iteration

A Simple example

A Simple example showing the implementation of k-means algorithm
(using K=2)

Individual	Variable 1	Variable 2
1	1.0	1.0
2	1.5	2.0
3	3.0	4.0
4	5.0	7.0
5	3.5	5.0
6	4.5	5.0
7	3.5	4.5

A Simple example

Step 1:

Initialization: Randomly we choose following two centroids ($k=2$) for two clusters.

In this case the 2 centroid are: $m1=(1.0,1.0)$ and $m2=(5.0,7.0)$.

Individual	Variable 1	Variable 2
1	1.0	1.0
2	1.5	2.0
3	3.0	4.0
4	5.0	7.0
5	3.5	5.0
6	4.5	5.0
7	3.5	4.5

	Individual	Mean Vector
Group 1	1	(1.0, 1.0)
Group 2	4	(5.0, 7.0)

A Simple example

Step 2:

Thus, we obtain two clusters containing:
 {1,2,3} and {4,5,6,7}.

Using Manhattan Distance

Their new centroids are:

	C1=(1.0,1.0)	C2=(5.0,7.0)	
(1.0,1.0)	1-1 + 1-1 =0	1-5 + 1-7 =10	C1
(1.5,2.0)	0.5+1=1.5	8.5	C1
(3.0,4.0)	2.0+3.0=5.0	5	C1
(5.0,7.0)	10	0	C2
(3.5,5.0)	6.5	3.5	C2
(4.5,5.0)	7	2.5	C2
(3.5,4.5)	6	4	C2

$$m_1 = \left(\frac{1}{3}(1.0 + 1.5 + 3.0), \frac{1}{3}(1.0 + 2.0 + 4.0) \right) = (1.83, 2.33)$$

$$m_2 = \left(\frac{1}{4}(5.0 + 3.5 + 4.5 + 3.5), \frac{1}{4}(7.0 + 5.0 + 5.0 + 4.5) \right) \\ = (4.12, 5.38)$$

A Simple example

Step 3:

Now using these centroids we compute the Euclidean distance of each object, as shown in table.

Therefore, the new clusters are:
 {1,2} and {3,4,5,6,7}

Next centroids are:

$m_1 = (1.25, 1.5)$ and $m_2 = (3.9, 5.1)$

	newC1(1.83,2.33)	newC2(4.12,5.38)	
(1.0,1.0)	0.83+1.33=2.16	7.5	C1
(1.5,2.0)	0.66	6	C1
(3.0,4.0)	2.84	2.5	C2
(5.0,7.0)	8.34	2.5	C2
(3.5,5.0)	4.34	1	C2
(4.5,5.0)	5.34	0.76	C2
(3.5,4.5)	3.84	1.5	C2

A Simple example

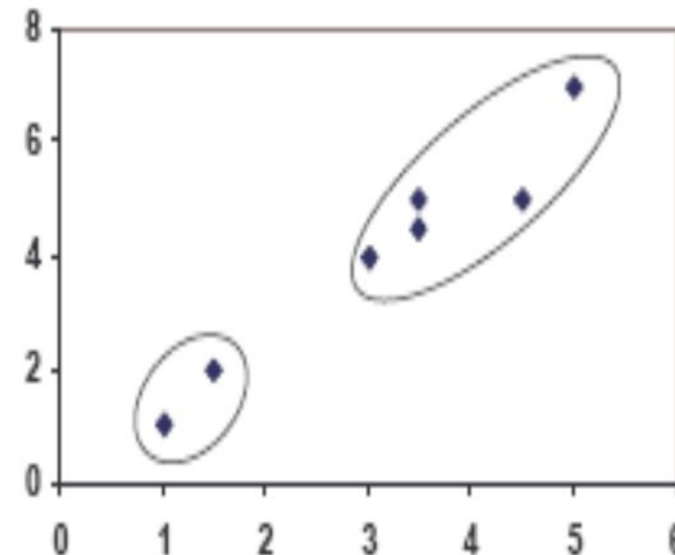
Step 4:

The clusters obtained are:
 $\{1,2\}$ and $\{3,4,5,6,7\}$

Therefore, there is no change in the cluster.

Thus, the algorithm comes to a halt here and final result consist of 2 clusters $\{1,2\}$ and $\{3,4,5,6,7\}$.

	newC1(1.25,1.5)	newC2(3.9,5.1)	
(1.0,1.0)	0.75	7	C1
(1.5,2.0)	0.75	5.5	C1
(3.0,4.0)	4.25	2	C2
(5.0,7.0)	9.75	3	C2
(3.5,5.0)	5.75	0.5	C2
(4.5,5.0)	6.75	0.7	C2
(3.5,4.5)	5.25	1	C2

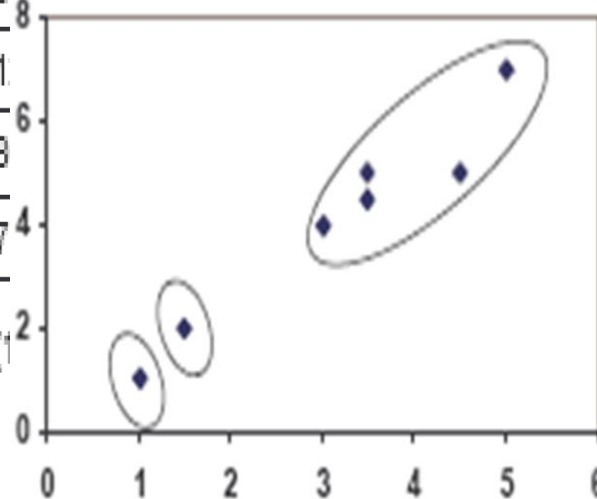


If $K=3$

Individual	$m_1 = 1$	$m_2 = 2$	$m_3 = 3$	cluster
1	0	1.11	3.81	1
2	1.12	0	2.5	2
3	3.81	2.5	0	3
4	7.21	6.10	3.81	3
5	4.72	3.81	1.1	3
6	5.31	4.24	1.8	3
7	4.30	3.20	0.7	3

clustering with initial centroids (1)

Step 1



Individual	m_1 (1.0, 1.0)	m_2 (1.5, 2.0)	m_3 (3.9, 5.1)	cluster
1	0	1.11	5.02	1
2	1.12	0	3.92	2
3	3.81	2.5	1.42	3
4	7.21	6.10	2.20	3
5	4.72	3.81	0.41	3
6	5.31	4.24	0.81	3
7	4.30	3.20	0.72	3

Step 2

K-Means Clustering

We have 4 medicines as our training data points object and each medicine has 2 attributes. Each attribute represents coordinate of the object. We have to determine which medicines belong to cluster 1 and which medicines belong to the other cluster.

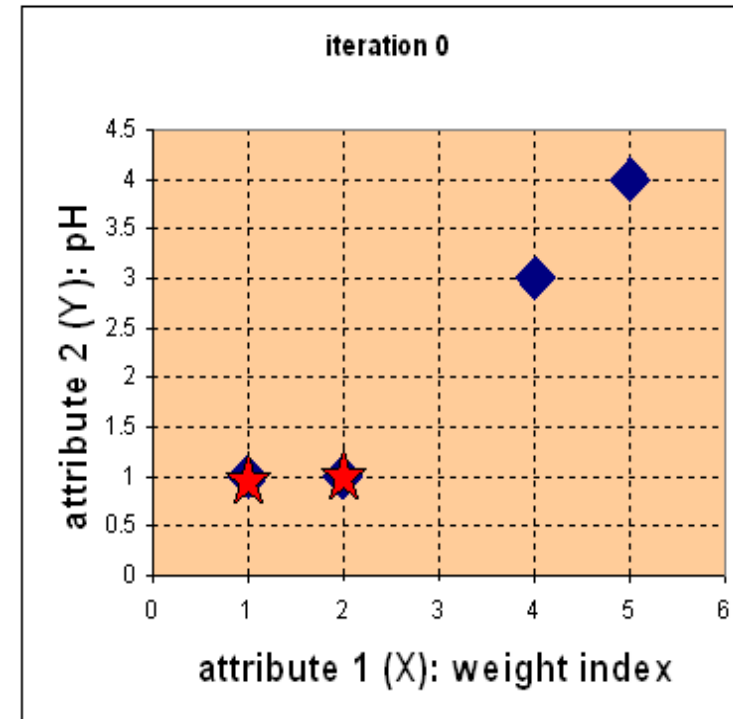
Object	Attribute1 (X): weight index	Attribute 2 (Y): pH
Medicine A	1	1
Medicine B	2	1
Medicine C	4	3
Medicine D	5	4

A Simple example - 2

Step 1:

Initial value of centroids : Suppose we use medicine A and medicine B as the first centroids.

Let c_1 and c_2 denote the coordinate of the centroids, then $c_1=(1,1)$ and $c_2=(2,1)$



A Simple example - 2

Step 1:

Objects-Centroids distance : we calculate the distance between cluster centroid to each object.

Let us use Euclidean distance, then we have distance matrix at iteration 0 is

$$D^0 = \begin{bmatrix} 0 & 1 & 3.61 & 5 \\ 1 & 0 & 2.83 & 4.24 \end{bmatrix} \quad \begin{matrix} c_1 = (1,1) \text{ group - 1} \\ c_2 = (2,1) \text{ group - 2} \end{matrix}$$

	A	B	C	D	
	1	2	4	5	X
	1	1	3	4	Y

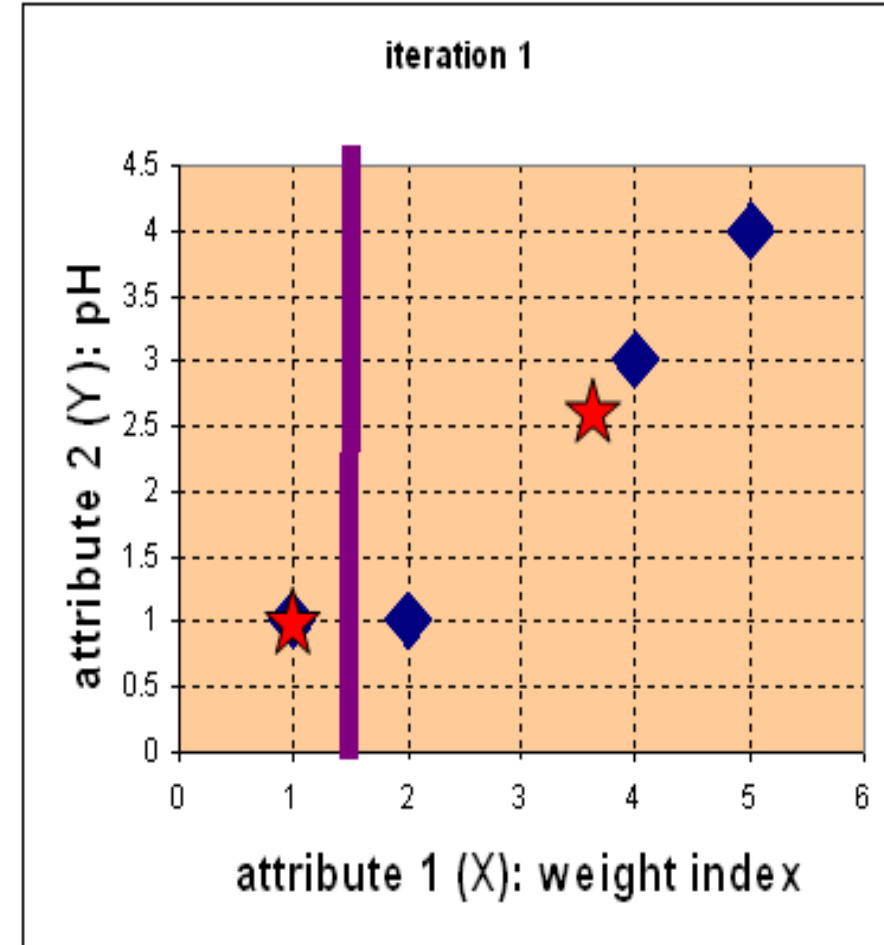
- Each column in the distance matrix symbolizes the object.
- The first row of the distance matrix corresponds to the distance of each object to the first centroid and the second row is the distance of each object to the second centroid.
- For example, distance from medicine C = (4, 3) to the first centroid $c_1 = (1,1)$ is , $\sqrt{(4-1)^2 + (3-1)^2} = 3.61$ and
- Its distance to the second centroid is , $c_2 = (2,1)$ is $\sqrt{(4-2)^2 + (3-1)^2} = 2.83$

A Simple example - 2

- **Step 2:**
- **Objects clustering** : We assign each object based on the minimum distance.
- Medicine A is assigned to group 1, medicine B to group 2, medicine C to group 2 and medicine D to group 2.
- The elements of Group matrix below is 1 if and only if the object is assigned to that group.

$$\mathbf{G}^0 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \end{bmatrix} \begin{matrix} \text{group - 1} \\ \text{group - 2} \end{matrix}$$

A B C D



A Simple example - 2

- **Iteration-1, Objects-Centroids distances :** The next step is to compute the distance of all objects to the new centroids.
- Similar to step 2, we have distance matrix at iteration 1 is

$$D^1 = \begin{bmatrix} 0 & 1 & 3.61 & 5 \\ 3.14 & 2.36 & 0.47 & 1.89 \end{bmatrix} \quad \begin{array}{l} \mathbf{c}_1 = (1,1) \text{ group-1} \\ \mathbf{c}_2 = (\frac{11}{3}, \frac{8}{3}) \text{ group-2} \end{array}$$

<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	
1	2	4	5	<i>X</i>
1	1	3	4	<i>Y</i>

A Simple example - 2

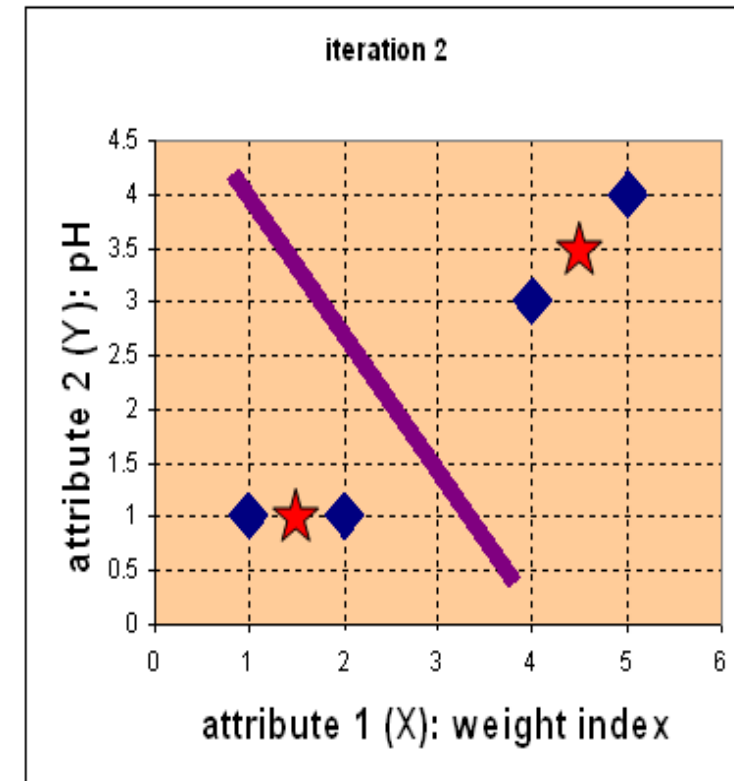
- **Iteration-1, Objects clustering:** Based on the new distance matrix, we move the medicine B to Group 1 while all the other objects remain. The Group matrix is shown below

$$\mathbf{G}^1 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \quad \begin{matrix} \text{group - 1} \\ \text{group - 2} \end{matrix}$$

$A \quad B \quad C \quad D$

- **Iteration 2, determine centroids:** Now we repeat step 4 to calculate the new centroids coordinate based on the clustering of previous iteration. Group1 and group 2 both has two members, thus the new centroids are

$$\text{and } \mathbf{c}_2 = \left(\frac{4+5}{2}, \frac{3+4}{2} \right) = (4\frac{1}{2}, 3\frac{1}{2}) \quad \mathbf{c}_1 = \left(\frac{1+2}{2}, \frac{1+1}{2} \right) = (1\frac{1}{2}, 1)$$



A Simple example - 2

- Iteration-2,
- Objects-Centroids distances : Repeat step 2 again, we have new distance matrix at iteration 2 as

$$\mathbf{D}^2 = \begin{bmatrix} 0.5 & 0.5 & 3.20 & 4.61 \\ 4.30 & 3.54 & 0.71 & 0.71 \end{bmatrix} \quad \begin{array}{l} \mathbf{c}_1 = (1\frac{1}{2}, 1) \text{ group - 1} \\ \mathbf{c}_2 = (4\frac{1}{2}, 3\frac{1}{2}) \text{ group - 2} \end{array}$$

	A	B	C	D	
	1	2	4	5	X
	1	1	3	4	Y

A Simple example - 2

- **Iteration-2,** Objects clustering: Again, we assign each object based on the minimum distance.
- We obtain result that $\mathbf{G}^2 = \mathbf{G}^1$. Comparing the grouping of last iteration and this iteration reveals that the objects does not move group anymore.
- Thus, the computation of the k-mean clustering has reached its stability and no more iteration is needed..

$$\mathbf{G}^2 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \begin{matrix} \text{group} - 1 \\ \text{group} - 2 \end{matrix}$$

$A \quad B \quad C \quad D$

A Simple example - 2

- We get the final grouping as the results as:

<u>Object</u>	<u>Feature1(X): weight index</u>	<u>Feature2 (Y): pH</u>	<u>Group (result)</u>
Medicine A	1	1	1
Medicine B	2	1	1
Medicine C	4	3	2
Medicine D	5	4	2

Applications of K-Mean Clustering

- It is relatively efficient and fast. It computes result at $O(tkn)$, where n is number of objects or points, k is number of clusters and t is number of iterations.
- k-means clustering can be applied to machine learning or data mining
- Used on acoustic data in speech understanding to convert waveforms into one of k categories (known as Vector Quantization or Image Segmentation).
- Also used for choosing color palettes on old fashioned graphical display devices and Image Quantization.

Weaknesses of K-Mean Clustering

- When the numbers of data are not so many, initial grouping will determine the cluster significantly.
- The number of cluster, K , must be determined before hand. Its disadvantage is that it does not yield the same result with each run, since the resulting clusters depend on the initial random assignments.
- We never know the real cluster, using the same data, because if it is inputted in a different order it may produce different cluster if the number of data is few.
- It is sensitive to initial condition. Different initial condition may produce different result of cluster. The algorithm may be trapped in the local optimum.

Advantages & Disadvantage of K means-Clustering

- Advantages
 - K-means is relatively scalable and efficient in processing large data sets
 - The computational complexity of the algorithm is $O(nkt)$
 - n : the total number of objects
 - k : the number of clusters
 - t : the number of iterations
 - Normally: $k \ll n$ and $t \ll n$
- Disadvantage
 - Can be applied only when the mean of a cluster is defined
 - Users need to specify k
 - K-means is not suitable for discovering clusters with non convex
 - Shapes or clusters of very different size
 - It is sensitive to noise and outlier data points

References

- *J. A. Hartigan (1975) "Clustering Algorithms". Wiley.*
- *J. A. Hartigan and M. A. Wong (1979) "A K-Means Clustering Algorithm", Applied Statistics, Vol. 28, No. 1, p100-108.*
- *www.wikipedia.com*



Dr. Vishwanath Karad

**MIT WORLD PEACE
UNIVERSITY** | PUNE

TECHNOLOGY, RESEARCH, SOCIAL INNOVATION & PARTNERSHIPS

Thank You