

## MODULE - 3.

### 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.

Classification - A two-step-process

- 1) Model Construction: describing a set of  
predefined 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 the training set.
- \* The model is represented as classification rules, decision trees or mathematical formulae.

2) Model Usage: for classifying future or unknown objects.

\* Estimate accuracy of the model.

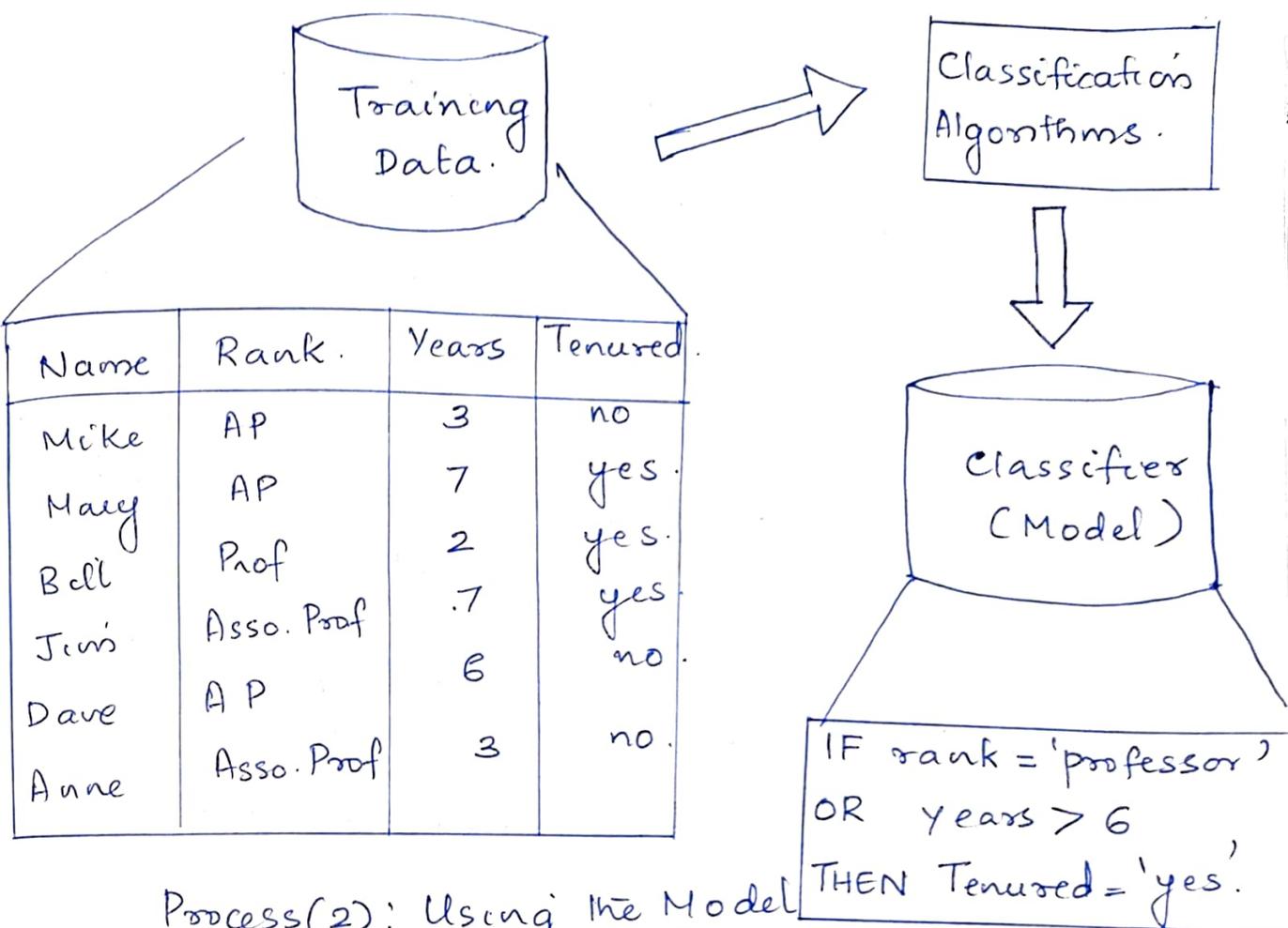
→ The known label of a test sample is compared with the classified result from the model.

→ Accuracy rate is the percentage of test set samples that are correctly classified by the model.

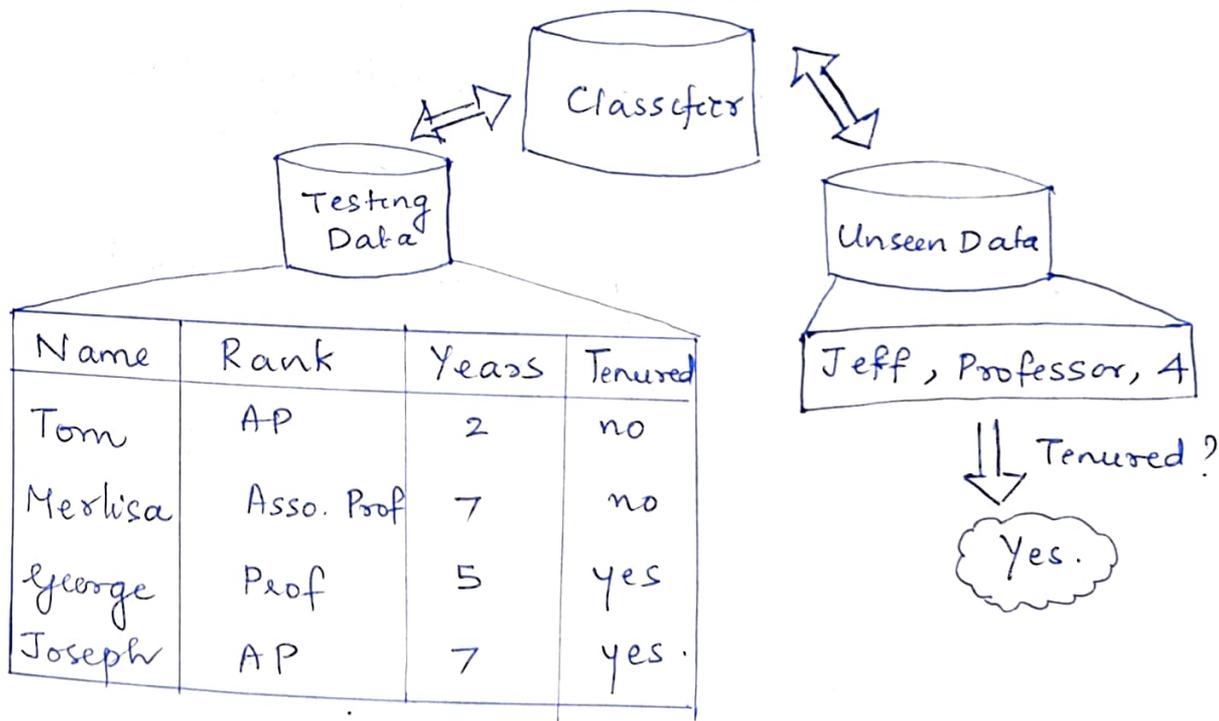
→ Test set is completely independent of training set, otherwise overfitting will occur.

\* If the accuracy is acceptable, use the model to classify data tuples whose class labels are not known.

eg: Process (1): Model Construction.



Process (2): Using the Model for Prediction.



## Supervised vs Unsupervised Learning.

### → Supervised Learning (classification)

Supervision: The training data (observations, measurements etc) are accompanied by labels indicating the class of observations.

New data is classified based on the training set.

### → Unsupervised Learning: (Clustering)

The class labels of training data are not known.

Given a set of observations, measurements etc with the aim of establishing the existence of classes or clusters in the data.

## Issues Regarding classification and Prediction.

### a) Preparing Data for classification and Prediction.

(i) Data cleaning: Preprocess the data in order to reduce noise and handle missing values.

ii) relevance analysis: Remove redundant and irrelevant attributes.

eg: Correlation analysis can be used to ~~find~~ identify whether 2 attributes are statistically related. Attribute Subset Selection can be used to find reduced set of attributes.

iii) data transformation and reduction:

- normalize data so that it lies within the range -1.0 to 1.0 or 0.0 to 1.0.
- data can be transformed by generalizing to higher level concepts.

eg: numeric values for attribute income can be generalized to discrete ranges such as low, medium and high.

b) Comparing Classification and Prediction Models.

(i) Accuracy:

accuracy of a classifier refers to the ability of a given classifier to correctly predict class label of new or previously unseen data.

accuracy of a predictor: refers to how well a predictor can guess the value of the predicted attribute for new or previously unseen data.

## ii) Speed.

- the computational costs involved in generating and using the classifier or predictor.

## iii) Robustness

- ability of the classifier or predictor to make correct predictions given noisy data or data with missing values.

## iv) Scalability.

- ability to construct the classifier or predictor efficiently given large amounts of data

## v) Interpretability.

- level of understanding & insight provided by the model.

## Decision Tree Induction.

- Decision Tree Induction is the learning of decision trees from class labelled training tuples

- A decision tree is a flowchart like structure

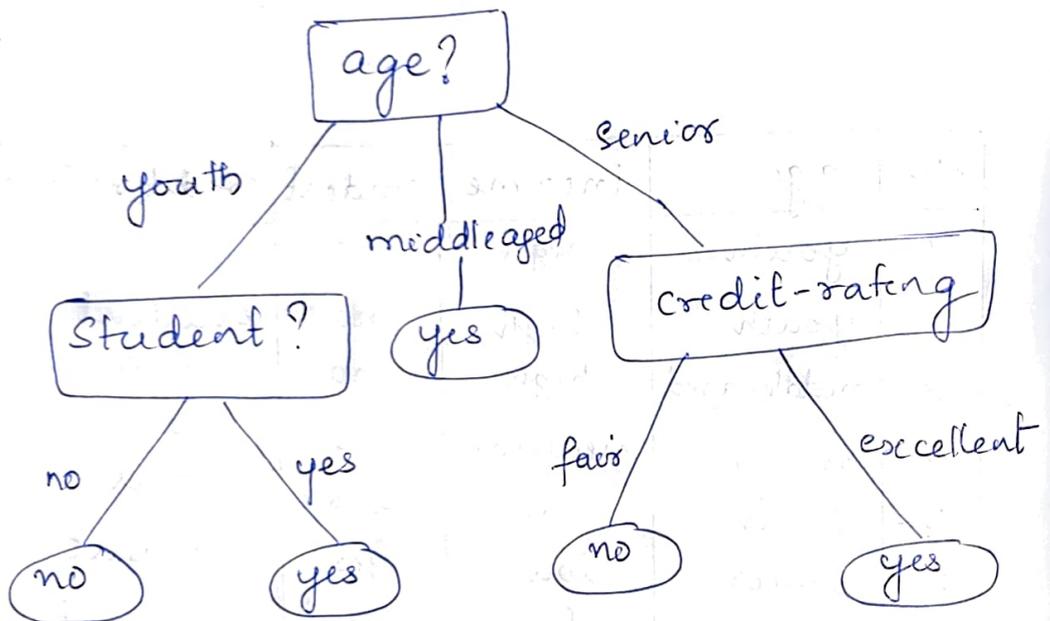
→ where each internal node denotes a test on an attribute

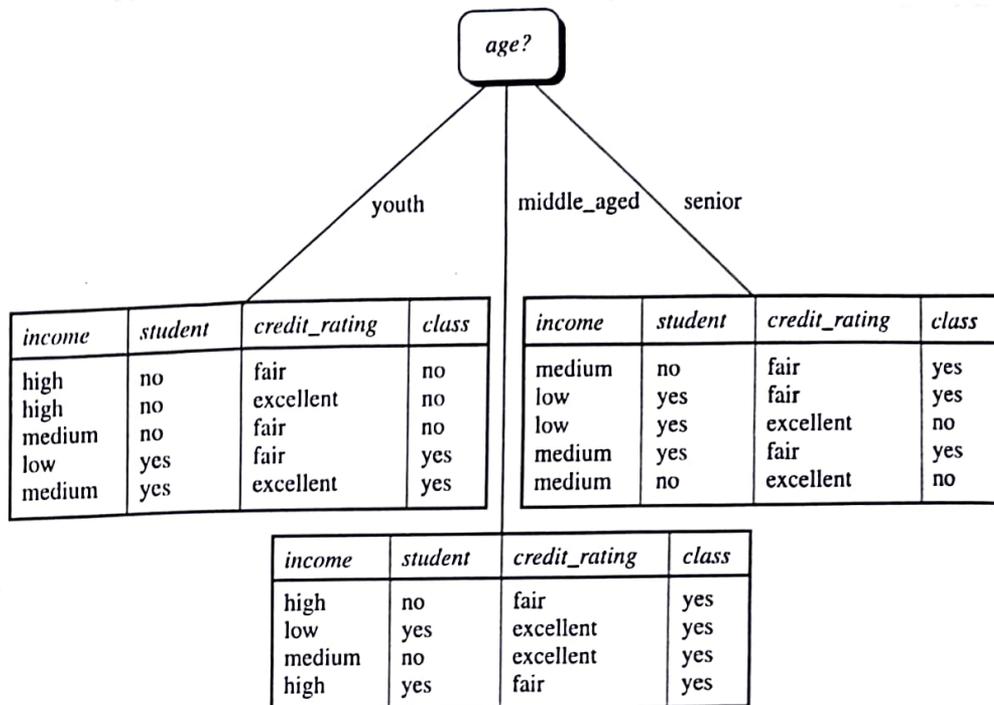
→ each branch represents an outcome of the test

→ each leaf node holds a class label.

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	no
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.

A decision tree for the concept `buys_computer`, indicating whether a customer at ACElectronics is likely to purchase a computer. Each internal (non leaf node) represents a test on an attribute. Each leaf node represents a class (either `buys_computer = yes` or `buys_computer = no`)





The attribute *age* has the highest information gain and therefore becomes the splitting attribute at the root node of the decision tree. Branches are grown for each outcome of *age*. The tuples are shown partitioned accordingly.

## Decision Tree - ID3.

- ID3 is a basic algorithm for learning Decision trees.

- Given a training set of samples, the algorithms for building decision trees performs search in the space of decision trees

- The construction of the tree is top-down and the approach is greedy.

- "Which attribute should be tested next?"  
Which attribute gives us more information?

- SELECT THE BEST ATTRIBUTE !!!!!

- A descendant node is then created for each possible value of this attribute and

Fundamental Question

examples are partitioned according to this value.

- The process is repeated for each successor node until all the examples are classified correctly or there are no attributes left.

Which attribute is the Best ???

- A statistical property called information gain measures how well a given attribute separates the training examples.

- Information gain uses the notion of entropy, commonly used in information theory.

- Information gain = expected reduction of entropy.

## I Attribute Selection Measure: Information Gain

- Select the attribute with the highest information gain.

- Let  $p_i$  be the probability that an arbitrary tuple in  $D$  belongs to class  $C_i$

used by  
ID3  
algorithm

estimated by  $P(C_i, D/D)$

What is Entropy???

Entropy is a measure of the randomness in the information being processed.

The higher the entropy, the harder is to draw any conclusion's from that information.

consider a segment 'S' of a dataset having 'c' number of class labels.

Let  $P_i$  be the proportion of examples in S having  $i^{\text{th}}$  class label.

Entropy is defined as

$$\text{Entropy}(S) = \sum_{i=1}^c -P_i \log_2(P_i)$$

Information Gain.

Information gain tells us how important a given attribute of the feature vector is.

- This value is used to decide the ordering of attributes in the nodes of a decision tree.

(S or D)  
S and D  
are same.

Let S be set of examples, A be a feature/attribute.

S<sub>v</sub> be subset of S with A = v and  
Values(A) be the set of all possible  
values of A.

Information Gain of an attribute A  
relative to S, denoted by Gain(S, A) is

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \times \text{Entropy}(S_v)$$

### ID3 Algorithm.

ID3(X, T, Attrs)

X: training examples  
T: target attribute  
Attrs: other attributes,  
Initially all attributes

Create Root node.

If all x's are +  
return Root with class +

If all x's are -  
return Root with class -

If Attrs is empty  
return Root with class most common value of  
T in X

Else

$A \leftarrow$  best attribute

decision attribute for Root  $\leftarrow A$

For each possible value  $v_i$  of  $A$ :

- add a new branch below Root, for test  $A = v_i$

-  $X_i \leftarrow$  subset of  $X$  with  $A = v_i$

- If  $X_i$  is empty then

add a new leaf with class the most common value of  $T$  in  $X$

else

add the subtree generated by

$ID3(X_i, T, Attrs - \{A\})$

return Root.

Class P: buys\_computer = "yes"

Class N: buys\_computer = "no"

(or)  $Info(S) = Entropy(S) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = \underline{\underline{0.940 \text{ bits}}}$

~~Entropy(age) =~~

$$Info_{age}(S) = \frac{5}{14} \times \left( -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) +$$
$$\frac{4}{14} \times \left( -\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} \right) +$$
$$\frac{5}{14} \times \left( -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right)$$

= 0.694 bits

$$\begin{aligned} \text{Gain}(S, \text{age}) &= \text{Info}(D) - \text{Info}_{\text{age}}(D) \\ &= 0.940 - 0.694 \\ &= \underline{\underline{0.246 \text{ bits}}} \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{income}}(S) &= \frac{4}{14} \left( -\frac{2}{4} \log_2 \left( \frac{2}{4} \right) - \frac{2}{4} \log_2 \left( \frac{2}{4} \right) \right) + \\ &\quad \frac{6}{14} \left( -\frac{4}{6} \log_2 \left( \frac{4}{6} \right) - \frac{2}{6} \log_2 \left( \frac{2}{6} \right) \right) + \\ &\quad \frac{4}{14} \left( -\frac{3}{4} \log_2 \left( \frac{3}{4} \right) - \frac{1}{4} \log_2 \left( \frac{1}{4} \right) \right) \\ &= 0.9097 \end{aligned}$$

$$\begin{aligned} \text{Gain}(S, \text{income}) &= 0.9403 - 0.9097 \\ &= \underline{\underline{0.03 \text{ bits}}} \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{student}}(S) &= \frac{7}{14} \left( -\frac{6}{7} \log_2 \left( \frac{6}{7} \right) - \frac{1}{7} \log_2 \left( \frac{1}{7} \right) \right) + \\ &\quad \frac{7}{14} \left( -\frac{3}{7} \log_2 \left( \frac{3}{7} \right) - \frac{4}{7} \log_2 \left( \frac{4}{7} \right) \right) \\ &= 0.7878 \end{aligned}$$

$$\begin{aligned} \text{Gain}(S, \text{student}) &= 0.9403 - 0.7878 \\ &= \underline{\underline{0.152 \text{ bits}}} \end{aligned}$$

Info credit-rating

$$I(S) = \frac{6}{14} \left( -\frac{3}{6} \log_2 \left( \frac{3}{6} \right) - \frac{3}{6} \log_2 \left( \frac{3}{6} \right) \right) + \frac{8}{14} \left( -\frac{6}{8} \log_2 \left( \frac{6}{8} \right) - \frac{2}{8} \log_2 \left( \frac{2}{8} \right) \right)$$

$$= 0.892$$

$$\text{Gain}(S, \text{credit-rating}) = 0.9403 - 0.892 = \underline{\underline{0.0483}}$$

Here the attribute with maximum Information gain is age with a value of 0.246 bits. So it will be chosen as the root of decision-tree.

The process is repeated for the rest of the attributes.

## II Attribute Selection Measure: Gain Ratio

used by C4.5

- Information gain measure is biased towards attributes with large number of values.
- C4.5 (a successor of ID3 algorithm) uses gain ratio to overcome the problem (normalization to information gain)

$$\text{SplitInfo}_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \log_2 \left( \frac{|D_j|}{|D|} \right)$$

$$\text{GainRatio}(A) = \frac{\text{Gain}(A)}{\text{SplitInfo}_A(D)}$$

The attribute with maximum gain ratio is selected as the splitting attribute.

(eg:-) To find the gain-Ratio of the attribute 'income'

$$\text{Gain}(S, \text{income}) = 0.029$$

$$\begin{aligned} \text{Split-Info}_{\text{income}}(D) &= -\frac{4}{14} \log_2 \left( \frac{4}{14} \right) - \\ &\quad -\frac{6}{14} \log_2 \left( \frac{6}{14} \right) - \\ &\quad -\frac{4}{14} \log_2 \left( \frac{4}{14} \right) \\ &= 0.926 \end{aligned}$$

$$\begin{aligned} \therefore \text{GainRatio}(\text{income}) &= \frac{\text{Gain}(D, \text{income})}{\text{Split-Info}_{\text{income}}(D)} \\ &= \frac{0.029}{0.926} = \underline{\underline{0.031}} \end{aligned}$$

### III Attribute Selection Measure: Gini Index.

used by  
CART

If a dataset 'D' contains examples from 'n' classes, gini index is defined

$$\text{as } \boxed{\text{gini}(D) = 1 - \sum_{j=1}^n p_j^2}$$

where  $p_j$  is the relative frequency of class j in D.

$$\boxed{\text{gini}_A(D) = \frac{|D_1|}{|D|} \text{gini}(D_1) + \frac{|D_2|}{|D|} \text{gini}(D_2)}$$

• Reduction in Impurity is defined as

$$\boxed{\Delta \text{gini}(A) = \text{gini}(D) - \text{gini}_A(D)}$$

- The attribute that maximizes the reduction in impurity (minimum Gini Index) is selected as the splitting attribute.

eg:  $\text{gini}(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = \underline{\underline{0.459}}$

- Now we need to find the gini index for each attribute.

eg: we start with "income"

consider the subset

$$S_1 = \{\text{low, medium}\}, \{\text{high}\}$$

Gini<sub>income</sub>  $\in S_1$

$$= \frac{10}{14} \text{Gini}(D_1) + \frac{4}{14} \text{Gini}(D_2)$$

$$= \frac{10}{14} \left( 1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2 \right) + \frac{4}{14} \left( 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 \right)$$

$$= \underline{\underline{0.443}}$$

$$S_2 = \{\text{low, high}\}, \{\text{medium}\}$$

Gini<sub>income</sub>  $\in S_2$

$$= \frac{8}{14} \text{Gini}(D_1) + \frac{6}{14} \text{Gini}(D_2)$$

$$= \frac{8}{14} \left[ 1 - \left(\frac{5}{8}\right)^2 - \left(\frac{3}{8}\right)^2 \right] + \frac{6}{14} \left[ 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 \right]$$

$$= \underline{\underline{0.4581}}$$

$$S_3 = \{\text{high, medium}\}, \{\text{low}\}$$

Gini<sub>income</sub>  $\in S_3$

$$= \frac{10}{14} \left[ 1 - \left(\frac{6}{10}\right)^2 - \left(\frac{4}{10}\right)^2 \right] + \frac{4}{14} \left[ 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \right]$$

$$= \underline{\underline{0.4500}}$$

Take the min value = 0.443

$$\begin{aligned}\therefore \Delta \text{Gini}(\text{income}) &= 0.459 - 0.443 \\ &= \underline{\underline{0.016}}\end{aligned}$$

The attribute that provides the smallest  $\text{gini}_A(D)$  (or largest reduction in impurity) is chosen to split the node.

### Comparison of the various Attribute Selection Measures.

1. Information gain.
  - biased towards multivalued attributes
2. Gain Ratio
  - tends to prefer unbalanced splits in which one partition is much smaller than others.
3. Gini Index
  - biased to multivalued attributes
  - has difficulty when no. of classes is large
  - tends to favor tests that result in equal sized partitions and purity in both partitions.

# Bayesian Classification.

- Statistical classifier.
- predicts class membership probabilities i.e., the probability that a given tuple belongs to a particular class.
- It is based on Bayes' theorem.
- Naive Bayes classifier is found to be comparable in performance with decision tree and other neural network classifiers.

Assumption

Naive Bayes' classifiers assume that the effect of an attribute value on a given class is independent of the values of other attributes

- It simplifies the computation and the term "naive"

## Bayes' Theorem

- Let  $X$  be a data tuple.
- In Bayesian terms,  $X$  is considered "evidence".  $X$  is described by a set of " $n$ " ~~tuples~~ attributes

- H be some hypothesis such that data tuple x belongs to some class 'C'.

- For classification problems, we want to determine  $P(H/x)$  is the probability that hypothesis "H" holds given the "evidence" or observed data "x".

- we are looking for the probability that a tuple x belongs to a class 'C' given the attribute description of x.

$$P(H/x) = \frac{P(x/H) \cdot P(H)}{P(x)}$$

$P(H/x)$  - posterior probability of H conditioned on x.

$P(H)$  - prior probability <sup>of H.</sup> (independent of x).

$P(x/H)$  - posterior probability of x conditioned on H.

$P(x)$  - prior probability of x

eg: Suppose x is a 35 year old customer with an income of \$40,000. H is hypothesis that our customer will buy a computer.

$P(H/x)$  - probability that customer x will buy a computer given the customer age + income

$P(H)$  - probability that any given customer will buy a computer regardless of age, income or any other information.

$P(X|H)$  - probability that a customer  $X$  is 35 years old and earns \$40,000 given that the customer will buy a computer.

$P(X)$  - probability that a person is 35 years old and earns \$40,000. from the set of customers

### Naive Bayesian classification: Working

1. Let  $D$  be a training set of tuples with associated class labels. Each tuple is a  $n$ -dimensional attribute vector.

$$X = (x_1, x_2, \dots, x_n)$$

2. Suppose that there are 'm' classes,  $C_1, C_2, \dots, C_m$ . Given a tuple  $X$ , the classifier will predict that  $X$  belongs to the class having the highest posterior probability conditioned on  $X$ .

i.e. we say  $X$  belongs to a class  $C_i$  iff

$$P(C_i|X) > P(C_j|X) \text{ for } 1 \leq j \leq m, j \neq i$$

The class  $C_i$  for which  $P(C_i|x)$  is maximized is called the maximum posterior hypothesis.

By Bayes' theorem

$$P(C_i|x) = \frac{P(x|C_i) \cdot P(C_i)}{P(x)}$$

3.  $P(x)$  is a constant for all classes and so only  $P(x|C_i) \cdot P(C_i)$  need to be maximized

→ If the class prior probabilities i.e.  $P(C_i)$  are not known, it is commonly assumed as  $P(C_1) = P(C_2) \dots = P(C_m)$

→ otherwise

$$P(C_i) = \frac{|C_i, D|}{|D|}$$

$|C_i, D|$  - no. of training tuples of class  $C_i$  in  $D$ .

4. Involves computation of  $P(x|C_i)$ .

Naive algorithm - assumption of class conditional independence - values of attributes are conditionally independent of one another.

$$P(x|C_i) = \prod_{k=1}^n P(x_k|C_i)$$

$$= P(x_1|C_i) \times P(x_2|C_i) \times \dots \times P(x_n|C_i)$$

To compute  $P(x|C_i)$  we consider the following:

- a) If  $A_k$  is categorical, then  $P(x_k|C_i)$  is no: of tuples of class  $C_i$  in  $D$  having value  $x_k$  for  $A_k$ , divided by  $|C_i, D|$  no: of tuples of class  $C_i$  in  $D$ .
- b) If  $A_k$  is continuous-valued, it is assumed to have a gaussian distribution with a mean  $\mu$  and standard deviation  $\sigma$  by.

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

eg. AllElectronics Customer DataBase.

Suppose  $X = (\text{age} = \text{youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit-rating} = \text{fair})$  to classify, tuple

By Bayes' Theorem

$$P(C_i|X) = P(X|C_i) \cdot P(C_i)$$

$$C_i \begin{cases} \text{buys-computer} = \text{yes} \\ \text{buys-computer} = \text{no} \end{cases}$$

To find  $P(C_i)$

$$P(C_1) = P(\text{buys-computer} = \text{yes}) = \frac{9}{14} = \underline{\underline{0.643}}$$

$$P(C_2) = P(\text{buys-computer} = \text{no}) = \frac{5}{14} = \underline{\underline{0.357}}$$

To find  $P(X | C_i)$

$$P(X | C_1) = P(X | \text{buys-computer} = \text{yes})$$

$$= P(\text{age} = \text{youth} | \text{buys-computer} = \text{yes}) \times$$

$$P(\text{income} = \text{medium} | \text{buys-computer} = \text{yes}) \times$$

$$P(\text{student} = \text{yes} | \text{buys-computer} = \text{yes}) \times$$

$$P(\text{credit-rating} = \text{fair} | \text{buys-computer} = \text{yes})$$

$$= \frac{2}{9} \times \frac{4}{9} \times \frac{6}{9} \times \frac{5}{9} = \underline{\underline{0.044}}$$

$$P(X | C_2) = P(X | \text{buys-computer} = \text{no})$$

$$= P(\text{age} = \text{youth} | \text{buys-computer} = \text{no}) \times$$

$$P(\text{income} = \text{medium} | \text{buys-computer} = \text{no}) \times$$

$$P(\text{student} = \text{yes} | \text{buys-computer} = \text{no}) \times$$

$$P(\text{credit-rating} = \text{fair} | \text{buys-computer} = \text{no}).$$

$$= \frac{3}{5} \times \frac{2}{5} \times \frac{1}{5} \times \frac{2}{5} = \underline{\underline{0.192}}$$

$$\therefore P(x | C_1) \times P(C_1)$$

$$= P(x | \text{buys-computer} = \text{yes}) \cdot P(\text{buys-computer} = \text{yes})$$

$$= 0.044 \times 0.643 = \underline{\underline{0.028}}$$

$$P(x | C_2) \times P(C_2)$$

$$= P(x | \text{buys-computer} = \text{no}) \cdot P(\text{buys-computer} = \text{no})$$

$$= 0.0192 \times 0.357 = \underline{\underline{0.0068}}$$

So we conclude that  $x$  belongs to class

$C_1$  i.e.  $\text{buys-computer} = \text{yes}$

Since  $P(C_1 | x) > P(C_2 | x)$

---

→ Laplacian Correction:

we have  $P(x | C_i) = P(x_1 | C_i) \times P(x_2 | C_i) \dots$

If any one of  $P(x_k | C_i) = 0$  then  $P(x | C_i)$  will be 0.

There is a simple trick to avoid this problem. We assume that our training database  $D$  is so large that adding one to each count would only make

a negligible difference in the estimated probability value, yet could conveniently avoid the case of probability value of zero. This technique for probability estimation is called Laplacian Correction / Laplacian Estimator.

eg: Suppose that for class buys-computer = yes in some DB contains 1000 tuples.

we have 0 tuples with income = low, 990 tuples with income = medium and 10 tuples with income = high.

$$P(\text{income} = \text{low} / \text{buys-computer} = \text{yes}) = \frac{0}{1000} = 0$$

$$P(\text{income} = \text{medium} / \text{buys-computer} = \text{yes}) = \frac{990}{1000} = 0.990$$

$$P(\text{income} = \text{high} / \text{buys-computer} = \text{yes}) = \frac{10}{1000} = 0.010$$

Using Laplacian Correction for the 3 quantities, we pretend that we have 1 more tuple to each income-value pair

Then the probabilities become

$$\frac{1}{1003} = \underline{\underline{0.001}}, \quad \frac{991}{1003} = \underline{\underline{0.988}}$$

$$\frac{11}{1003} = \underline{\underline{0.011}}$$