

Classification - A two step process

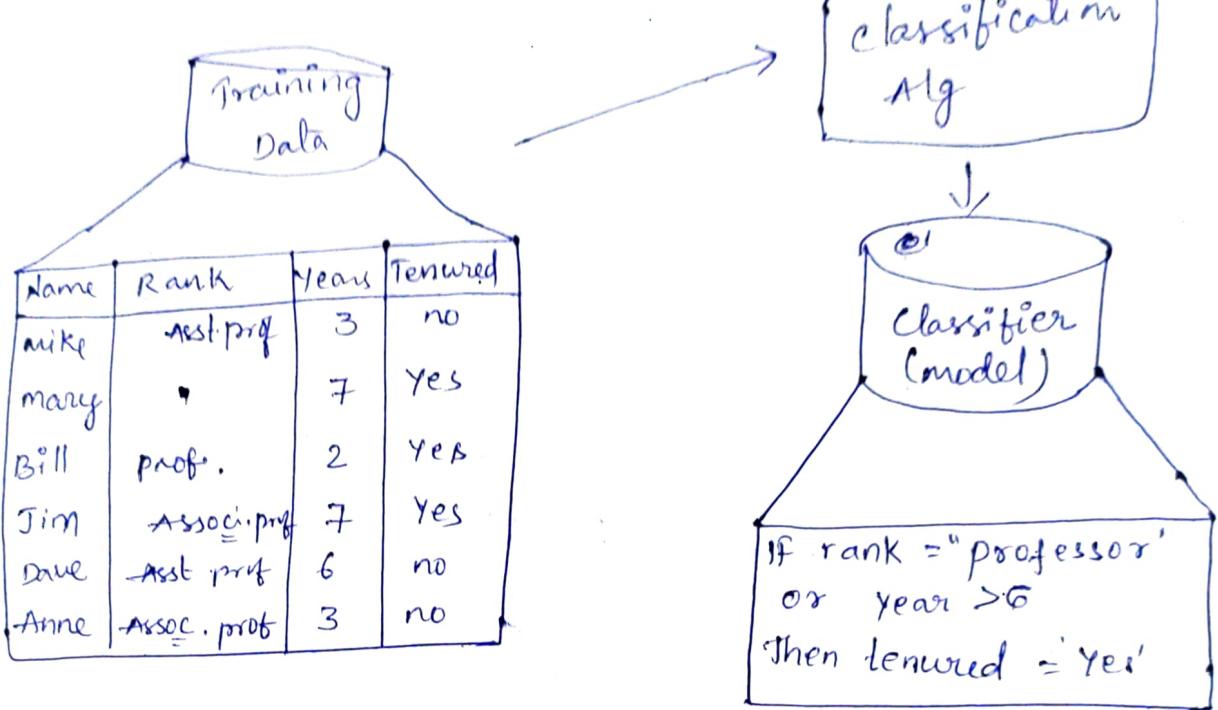
③ Model construction: describing a set of predetermined classes.

- Each tuple / sample is assumed to belong to a predefined class, as determined by the classlabelattribute.
- The set of tuples used for model construction is training set:
- The model is represented as classification rules, decision trees, or mathematical formulae.

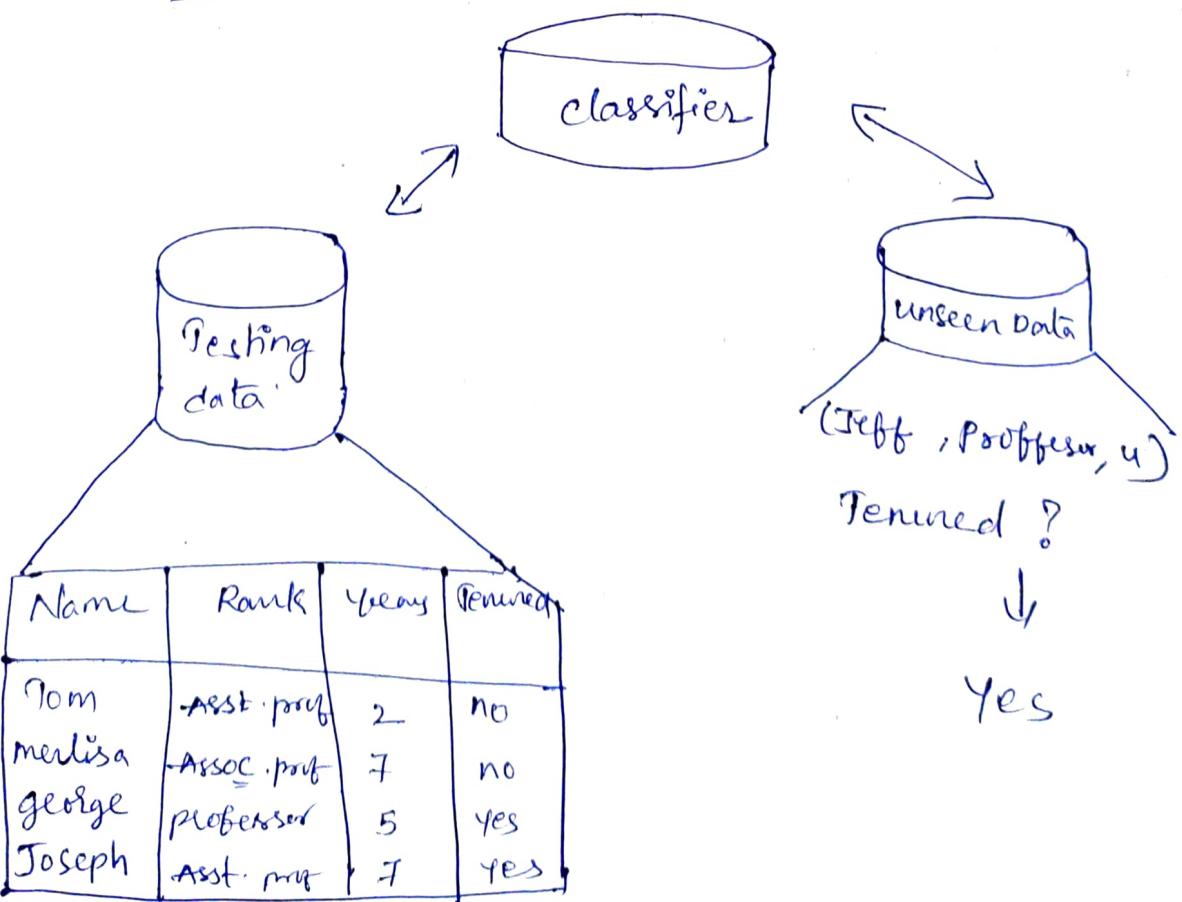
② Model usage: for classifying future or unknown objects.

- Estimate accuracy of the model.
 - ① The known label of the test sample is compared with the classified result from the model.
 - ② Accuracy rate is the percentage of test set samples that are correctly classified by the model.
- ③ 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.

Process ① : Model Construction

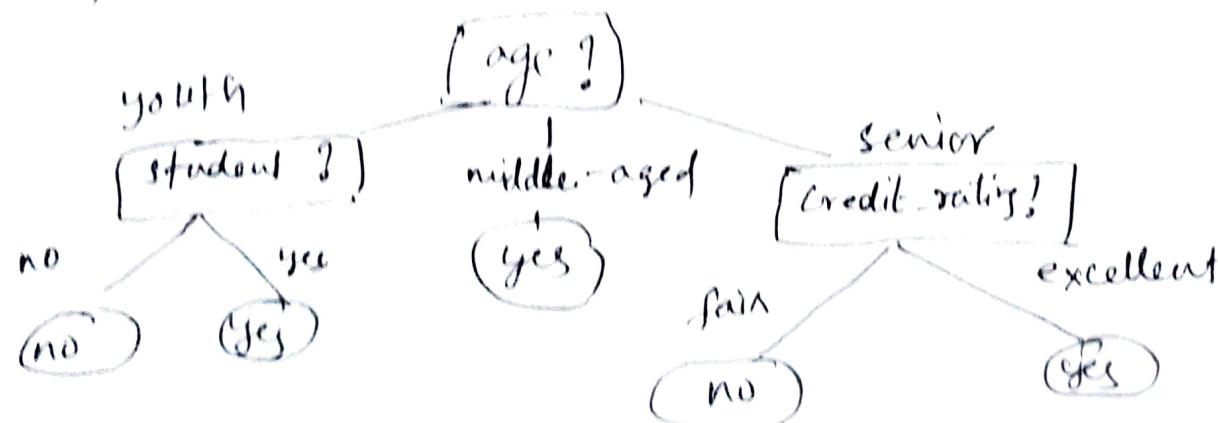


Process ② : Using the model in prediction



Classification by Decision Tree Induction

- Decision tree induction is the learning of decision trees from class-labeled training tuples.
- A decision tree is a flow-chart like tree structure, where each internal node (nonleaf node) denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node (or terminal node) holds a class label. The top-most node is also the root node.



Decision Tree Induction

Alg: Generate - decision - tree . Generate a decision tree from the training tuples of data partition D .

INPUT :

- Data partition, D, which is a set of training tuples and their associated class labels;
- attribute-list, the set of candidate attributes;
- Attribute-selection-method, a procedure to determine the splitting criterion that "best" partitions the data tuples into individual classes. This criterion consists of a splitting-attribute and possibly, either a split point or splitting subset .

OUTPUT : A decision tree

Method :

- ① create a node N;
- ② if tuples in D are all of the same class, C then
- ③ return N as a leaf node labeled with the class C;
- ④ if attribute-list is empty then
- ⑤ return N as a leaf node labeled with the majority class in D; // majority voting
- ⑥ apply Attribute-Selection-method (D, attribute-list) to find the "best" splitting-criterion;
- ⑦ label node N with splitting-criterion;
- ⑧ if splitting-attribute is discrete-valued and multiway splits allowed then // not restricted to binary tree
- ⑨ attribute-list \leftarrow attri-list - splitting-attribute; // remove splitting-attribute
- ⑩ for each outcome j of splitting-criterion.
// partition the tuples and grow subtrees for each partition

let D_j be the set of data tuples in D satisfying outcome j ; // a partition.
 if D_j is empty then
 attach a leaf labeled with the majority class in D to node N ,
 else attach the node returned by generate-decision-tree
 (D_j , attribute-list) to node N ;
 end for
 return N ;

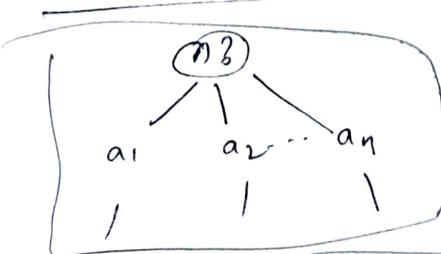
A discrete-valued
 A is continuous-valued
 A is discrete-valued and a binary tree

fig: Basic Alg for inducing a decision tree from training tuples.

→ the alg is called with 3 parameters:
 D , attribute-list and Attribute-selection-method
 where D — data partition, initially, it is the complete set of training tuples and their associated class labels.
 → The parameter attribute-list is a list of attributes describing the tuples.
 → Attribute-selection-method specifies a heuristic procedure for selecting the attribute that "best" discriminates the given tuples according to class.

This procedure employs an attribute selection measure, such as information gain or the Gini Index.

partitioning scenarios



Examples



$A \leq \text{split-point}$ $A > \text{split-point}$

Splitting Criterion

Attribute selection Measures

→ An attribute selection measure is a heuristic for selecting the splitting criterion that "best" separates a given data partition, D , of class-labeled training tuples into individual classes.

Information gain

= Iterative Dichotomizer
→ ID₃ uses info gain as its attribute selection measure. This measure is based on pioneering work by Claude Shannon on info theory, which studied the value or "info content" of messages. Let node N rep. or hold the tuples of partition D . The attr. with the highest info gain is chosen as the splitting attr. for the node N . This attr. minimizes the info needed to classify the tuples in the resulting partitions and reflects the least randomness or "impurity" in these partitions.

→ Such an approach minimizes the expected no. of tests needed to classify a given tuple and guarantees that a simple (but not necessarily the simplest) tree is found.

The expected info needed to classify a tuple in D is given by

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

- where P_i is the probability that an arbitrary tuple in D belongs to class C_i and is estimated by $|C_i, D| / |D|$
- A log function to the base is used, because the information is encoded in bits.
- $\text{Info}(D)$ is just the average amount of information needed to identify the class label of a tuple in D .
- Note that, at this point, the info we have is based solely on the proportions of tuples of each class.
- $\text{Info}(D)$ is also known as the entropy of D .

Now, suppose we were to partition the tuples in D on some attribute A having V distinct values, $\{a_1, a_2 \dots a_V\}$ as observed from the training data

- If A is discrete-valued, these values correspond directly to the V outcomes of a test on A .

- Attribute A can be used to split D into v partitions or subsets, $\{D_1, D_2 \dots D_v\}$, where D_j contains those tuples in D that have outcome a_j of A.
- These partitions would correspond to the branches grown from node N.
- Ideally, we would like this partitioning to produce an exact classification of the tuples.
- That is, we would like for each partition to be pure.
- However, it is quite likely that the partitions will be impure (e.g. where a partition may contain a collection of tuples from diff classes rather than from a single class)
- How much more info would we still need (after the partitioning) in order to arrive at an exact classification?
- This amount is measured by

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

- The term $\frac{|D_j|}{|D|}$ acts as the weight of the j^{th} partition.
- $Info(D)$ is the expected $info_{=}$ required to classify a tuple from D based on the partitioning by A .
- Information gain is defined as the difference b/w the original $info_{=}$ requirement (i.e. based on just the proportion of classes) and the new requirement (i.e., obtained after partitioning on A).

→ That is

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

- $Gain(A)$ tells us how much would be gained by branching on A .

$$\text{Info}(D) = -\frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right) = 0.409 + 0.530$$

$\log_2 0.639 = -0.48$

$$\approx 0.940 \text{ bits}$$

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

Hence, the gain in information from such a partitioning would be

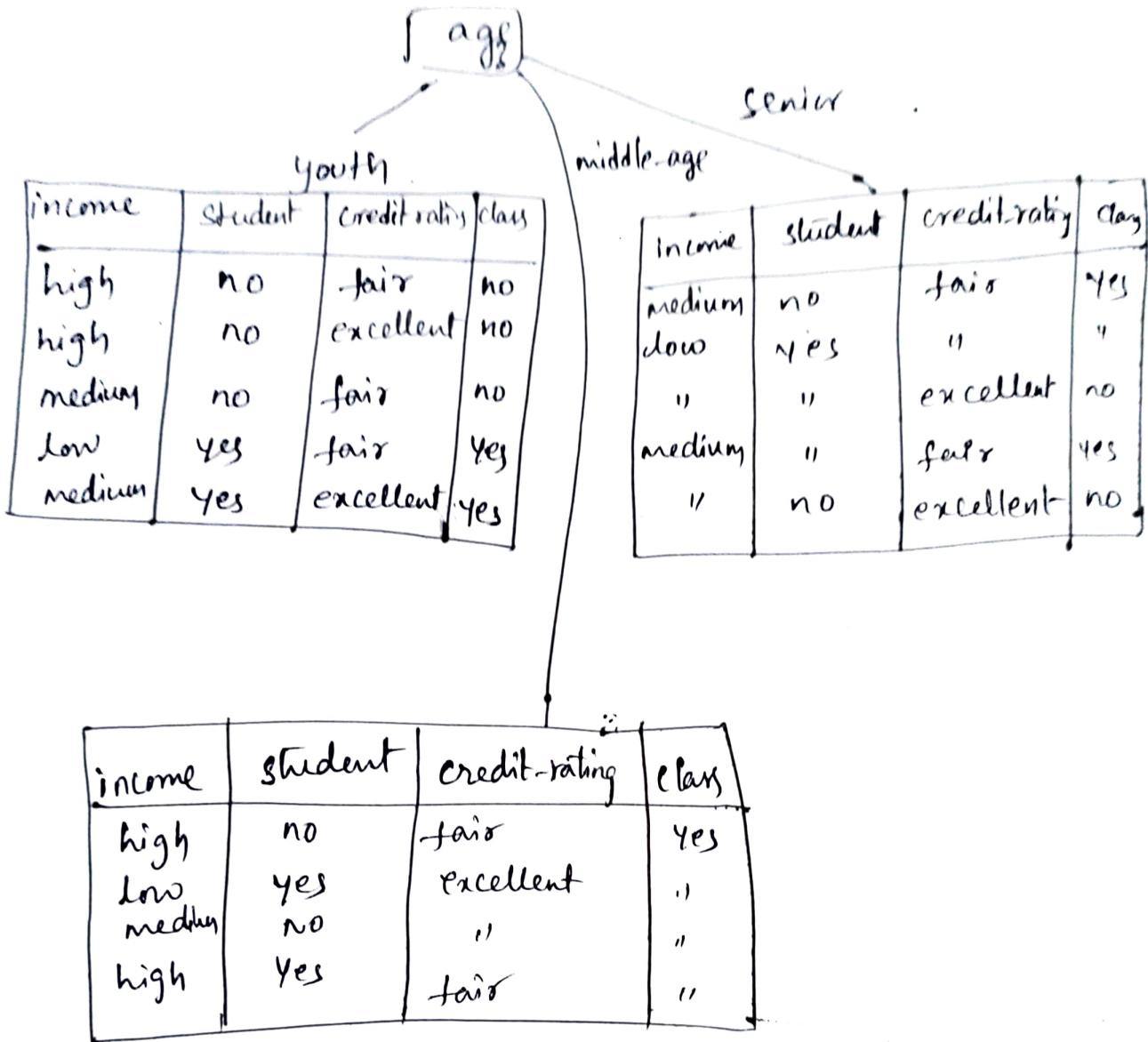
$$\begin{aligned}\text{Gain}(\text{age}) &= \text{Info}(D) - \text{Info}_{\text{A}}(D) \\ &= 0.940 - 0.694 = 0.246 \text{ bits}\end{aligned}$$

Similarly, we compute $\text{Gain}(\text{income}) = 0.029 \text{ bits}$

$\text{Gain}(\text{student}) = 0.151 \text{ bits}$

$\text{Gain}(\text{credit-rating}) = 0.048 \text{ bits}$

Be'coz age has the highest info gain among the attributes, it is selected as the splitting attribute.



entropy. $S = [9+, 5-] = 0.94$ / impurity

$$S_{youth} = [2+, 2-] = 1$$

$$S_{middle} = [4+, 2-] = 0.91$$

$$S_{senior} = [4+, 0] = 0$$
 - pure

purity

2+, 2-

Gain ratio

- The info gain measure is biased toward tests with many outcomes.
- That is, it prefers to select attributes having a large no. of values.
- For ex: consider an attribute that acts as a unique identifier, such as product-ID. A split on product-ID would result in a large no. of partitions (as many as there are values), each one containing just one tuple.
- Becoz each partition is pure, the info required to classify data set D based on this partitioning ~~or~~ would be $\boxed{\text{Info}_{\text{Product-ID}}(D) = 0}$
- Therefore, the info gained by partitioning on this attribute is maximal.
- C4.5, a successor of ID3, uses an extension to info gain known as gain ratio, which attempts to overcome this bias. It applies a kind of normalization to info gain using a "Split information" value defined analogously with $\text{Info}(D)$ as

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

... (1)

→ This value represents the potential info generated by splitting the training data set, D , into V partitions, corresponding to the V outcomes of a test on attribute A .

Note that, for each outcome, it considers the no. of tuples having that outcome with respect to the total no. of tuples in D .

→ It differs from info gain, which measures the info w.r.t. classification that is acquired based on the same partitioning.

The gain ratio is defined as

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

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

Ex: In above example , A test on income splits the data into three partitions , namely low , medium & high containing four , six and four tuples resp .

→ To compute the gain ratio of income , we first use eqn ① to obtain

$$\text{SplitInfo}_A(D) = -\frac{4}{14} \times \log_2 \left(\frac{4}{14} \right) - \frac{6}{14} \times \log_2 \left(\frac{6}{14} \right) - \frac{4}{14} \times \log_2 \left(\frac{4}{14} \right)$$

|| = 0.926.

we have $\underset{A}{\text{Gain(income)}} = 0.029$

$$\therefore \text{Gain Ratio(income)} = \frac{0.029}{0.926} = 0.031.$$

Gini Index

Classification & Regression
Trees

→ the Gini Index is used in CART . Using the the Gini Index measures the impurity of D , a data partition or set of training tuples , as

$$\text{Gini}(D) = 1 - \sum_{i=1}^m p_i^2$$

where p_i is the probability that a tuple in D belongs to class C_i and is estimated by $|C_i, D| / |D|$. The sum is computed over m classes

When considering a binary split, we compute a weighted sum of the impurity of each resulting partition.

Ex: if a binary split on A partitions D into D_1 and D_2 , given that partitioning is the gini index of D.

$$Gini_A(D) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2)$$

$$\Delta \text{Gini}(A) = \text{Gini}(D) - \text{Gini}_A(D) \quad (2)$$

Now let D be the training data of above table, where there are nine tuples belonging to the class $\text{buys-computer} = \text{yes}$ and the remaining five tuples belong to the class $\text{buys-computer} = \text{no}$. A (root) node N is created for the tuples in D . we use eqn(1) for GINI index to compute the impurity of D :

$$\text{Gini}(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.4159.$$

To find the splitting criterion for the tuples in D, we need to compute the gini index for each attribute.

$$\begin{aligned}
 & \text{Gini}_{\text{income}} \text{ of low, medium } (D) \\
 & = \frac{10}{14} \text{ Gini}(D_1) + \frac{4}{14} \text{ Gini}(D_2) \\
 & = \frac{10}{14} \left(1 - \left(\frac{6}{10}\right)^2 - \left(\frac{4}{10}\right)^2\right) + \frac{4}{14} \left(1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2\right) \\
 & \approx 0.450
 \end{aligned}$$

$$= \text{Gini}_{\text{Income}} \in \{\text{high}\} \quad (0)$$

- Select the best attribute which gives the minimum gini index overall.

- Bayesian classification
- Bayesian classifiers are statistical classifiers. They can predict class membership probabilities, such as the prob that a given tuple belongs to a particular class.
- Bayesian to be comparable in performance with decision tree and selected neural n/w classifier.
- also exhibited high accuracy & speed when applied to large db's.

- Naïve Bayesian classifiers independent of the values of the other attributes this assumption is called class conditional independence.

Bayes theorem

- Let 'x' be a data tuple .
- X is considered "evidence".
- Let H be some hypothesis, such as that the data tuple X belongs to a specified class C.
- $P(H/x)$ is the posterior prob of H conditioned on X.
- X is 35-yr old customer with an income of \$40,000.
- Suppose that H is the hypothesis that our customer will buy a computer.

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

Naïve Bayesian classification or simple Bayesian classifier

- ① let D be a training set of tuples and their associated class labels. $x = (x_1, x_2, \dots, x_n)$ A_1, A_2, \dots, A_n n attributes
 - ② m classes C_1, C_2, \dots, C_m , given that X belongs to the class having the highest posterior probability, the classifier will predict that X belongs to $P(C_i/x) > P(C_j/x)$ for $1 \leq j \leq m, j \neq i$.
 - ③ Thus we maximize $P(C_i/x)$ is called maximum posterior probability.
- $$P(C_i/x) = \frac{P(x/C_i) P(C_i)}{P(x)}$$
- As $P(x)$ is constant for all classes, only $P(x/C_i) P(C_i)$ need be maximized.

(4) In order to reduce computation in evaluating $P(X|c_i)$ the naive assumption of class conditional independence is made.

$$P(X|c_i) = \prod_{k=1}^n P(x_k|c_i)$$

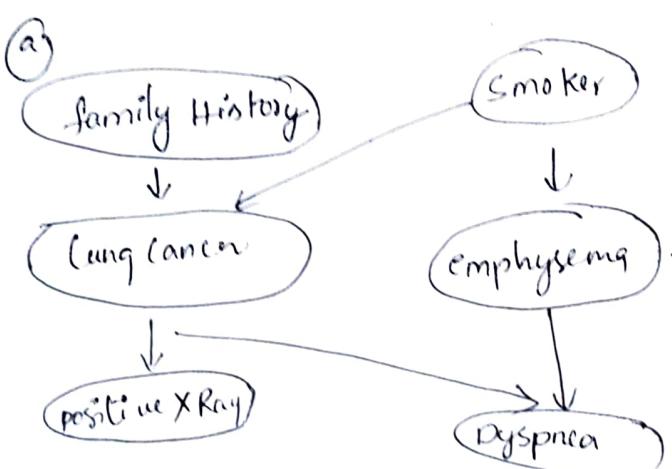
$$= P(x_1|c_i) \times P(x_2|c_i) \times \dots \times P(x_n|c_i)$$

(5) In order to predict the class label of x , $P(X|c_i) P(c_i)$ is evaluated for each class c_i .

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

(3) Bayesian Belief Network :- is defined by 2 components

- (1) a directed acyclic graph
- (2) a set of conditional probability table.



(b)

	FH, S	FH, ~S	~FH, S	~FH, ~S
LC	0.8	0.5	0.7	0.6
-LC	0.2	0.5	0.3	0.4

fig: the conditional probability table for the values of the variables

LungCancer (LC) showing each possible combination of the values of its parent nodes, FamilyHistory (FH) & Smoker (S).

(2) fig: A proposed causal model, represented by a directed acyclic graph.

Example : $X = (\text{age} = \text{youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit-rating} = \text{fair})$, predict the class label using naive bayesian classifier.

RID	age	income	student	credit-rating	class	Naive Bayesian classification
1	Youth	high	no	fair	no	
2	"	"	"	excellent	"	
3	middle-aged	"	"	fair	yes — (1)	
4	senior	medium	"	"	" — (2)	
5	"	low	yes	"	" — (3)	
6	"	"	"	Excellent	no	
7	middle-aged	"	"	"	yes — (4)	
8	Youth	medium	no	fair	no	
9	"	low	yes	"	yes — (5)	
10	Senior	medium	"	"	" — (6)	
11	Youth	"	"	Excellent	" — (7)	
12	middle-aged	"	no	"	" — (8)	
13	"	high	yes	fair	" — (9)	
14	Senior	medium	no	Excellent	no	

Prior probabilities

$$\bar{P}(\text{yes}) = 9/14 = 0.643$$

$$P(\text{no}) = 5/14 = 0.357$$

Probability of likelihood

$$P(\text{Youth} | \text{yes}) = 2/9 = 0.222$$

$$P(\text{Youth} | \text{no}) = 3/5 = 0.600$$

$$P(\text{credit-rating} = \text{fair} | \text{yes}) = 6/9 = 0.667$$

$$P(\text{credit-rating} = \text{fair} | \text{no}) = 2/5 = 0.400$$

$$P(\text{income} = \text{medium} | \text{yes}) = \frac{4}{9} = 0.444$$

$$P(\text{income} = \text{medium} | \text{no}) = \frac{2}{5} = 0.400$$

$$P(\text{student} = \text{yes} | \text{yes}) = \frac{6}{9} = 0.667$$

$$P(\text{student} = \text{no} | \text{no}) = 1/5 = 0.200$$

$x \Rightarrow (\text{age} = \text{youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit-rating} = \text{fair})$

$$\begin{aligned} & P(\text{yes} | \text{youth, medium, yes, fair}) \\ p &= \frac{[(\text{youth} | \text{yes}) * P(\text{medium} | \text{yes}) * P(\text{yes} | \text{yes}) * P(\text{fair} | \text{yes})]}{P(\text{yes})} \\ &= (0.222 * 0.644 * 0.667 * 0.667) * 0.643 \\ &= \underline{\underline{0.028}} \end{aligned}$$

$$\begin{aligned} & P(\text{no} | \text{youth, medium, yes, fair}) \\ &= \frac{[P(\text{youth} | \text{no}) * P(\text{medium} | \text{no}) * P(\text{yes} | \text{no}) * P(\text{fair} | \text{no})]}{P(\text{no})} \\ &= (0.600 * 0.600 * 0.200 * 0.400) * 0.357 \\ &= \underline{\underline{0.007}} \end{aligned}$$

\therefore The naive bayesian classifier predicts buys-computer = yes for tuple x .

Bayesian Belief Networks

- The naive Bayesian classifier makes the assumption of class conditional independence, that is, given the class label of a tuple, the values of the attributes are assumed to be conditionally independent of one another.
- This simplifies computation.
- Bayesian belief nets specify joint conditional probability distributions. They allow class conditional independencies to be defined b/w subsets of variables.
- They provide a graphical model of causal relationships, on which learning can be performed.
- Bayesian belief nets are also known as belief nets, Bayesian nets and probabilistic nets.
- A belief net is defined by two components
 - ① A directed acyclic graph.
 - ② A set of conditional probability tables.
- Each node in the directed acyclic graph represents a random variable.
- The variables may be discrete or continuous-valued.

- They may correspond to actual rotation attributes given in the data or to "hidden variables" believed to form a relationship
ex: in the case of medical data, a hidden variable may indicate a syndrome, representing a no. of symptoms that, together, characterize a specific disease).
- In fig allow a representation of causal knowledge.
Ex: having lung cancer is influenced by a person's family history of lung cancer, as well as whether or not the person is a smoker.
Note: that the variable Positive XRay is independent of whether the patient has a family history of lung cancer or is a smoker, given that we know the patient has lung cancer.
- In otherwords, once we know the outcome of the variable LungCancer, then the vars FamilyHistory & smoker do not provide any additional info regarding Positive XRay.
- The arcs also show that the var LungCancer is conditionally independent of Emphysema, given its parents, familyHistory & Smoker.

- A belief net has one conditional probability table (CPT) for each var. The CPT for a var y specifies the conditional distribution $P(y | \text{parents}(y))$, where parents(y) are the parents of y .
- Fig (b) shows a CPT for the variable LungCancer.
- The conditional probability for each known value of LungCancer is given for each possible combination of values of its parents.

→ For instance, from the upper leftmost and bottom-rightmost entries, resp., we see that-

$$P(\text{LungCancer} = \text{Yes} | \text{FamilyHistory} = \text{Yes}, \text{Smoker} = \text{Yes}) = 0.8$$

$$P(\text{LungCancer} = \text{no} | \text{FamilyHistory} = \text{no}, \text{Smoker} = \text{no}) = 0.9$$

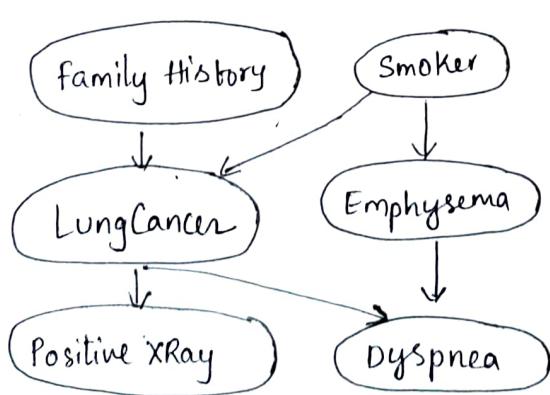


fig: A proposed causal model, represented by a directed acyclic graph.

	FH, S	FH, ~S	~FH, S	~FH, ~S
LC	0.8	0.5	0.7	0.1
-LC	0.2	0.5	0.3	0.9

fig: The conditional probability table for the values of the var LungCancer(LC) showing each possible combination of the values of its parent nodes, FamilyHistory(FH) and Smoker(S).