

Module - 4 (Classification Analysis) ①

Classification

classification is the process of finding a model that predicts the class label like categorical values.

- Example :-
- ① A bank loan officer needs analysis of her data to learn which loan applicants are safe (S) or risky.
 - ② A marketing manager at All-electronics needs data analysis to help whether customers with a given profile will buy a new computer. class labels like Yes (S) or No.
 - ③ A medical researcher wants to analyze breast cancer data to predict which one of three specific treatments a patient should receive. class labels like Treatment A, Treatment B, Treatment C for medical data.

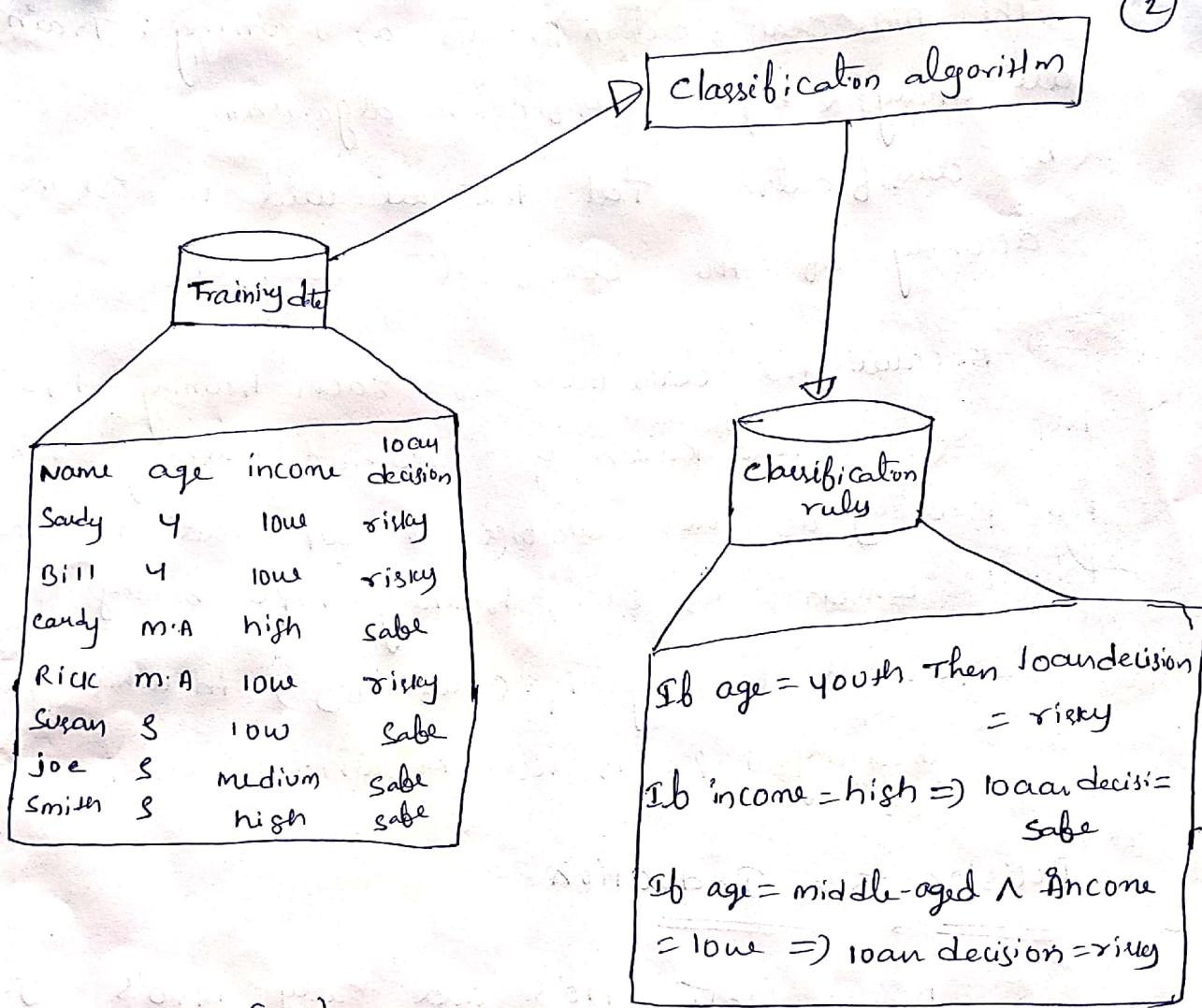
In each of these example Data Analysis task is classification, where a model (S) classifying is conducted to predict class label.

If suppose that the marketing manager wants to predict how much of a given customer will spend during a sale at All-electronics. The data Analysis task is numeric prediction. this is called Regression Analysis classification & Regression is two popular terms for prediction.

General Approach to classification

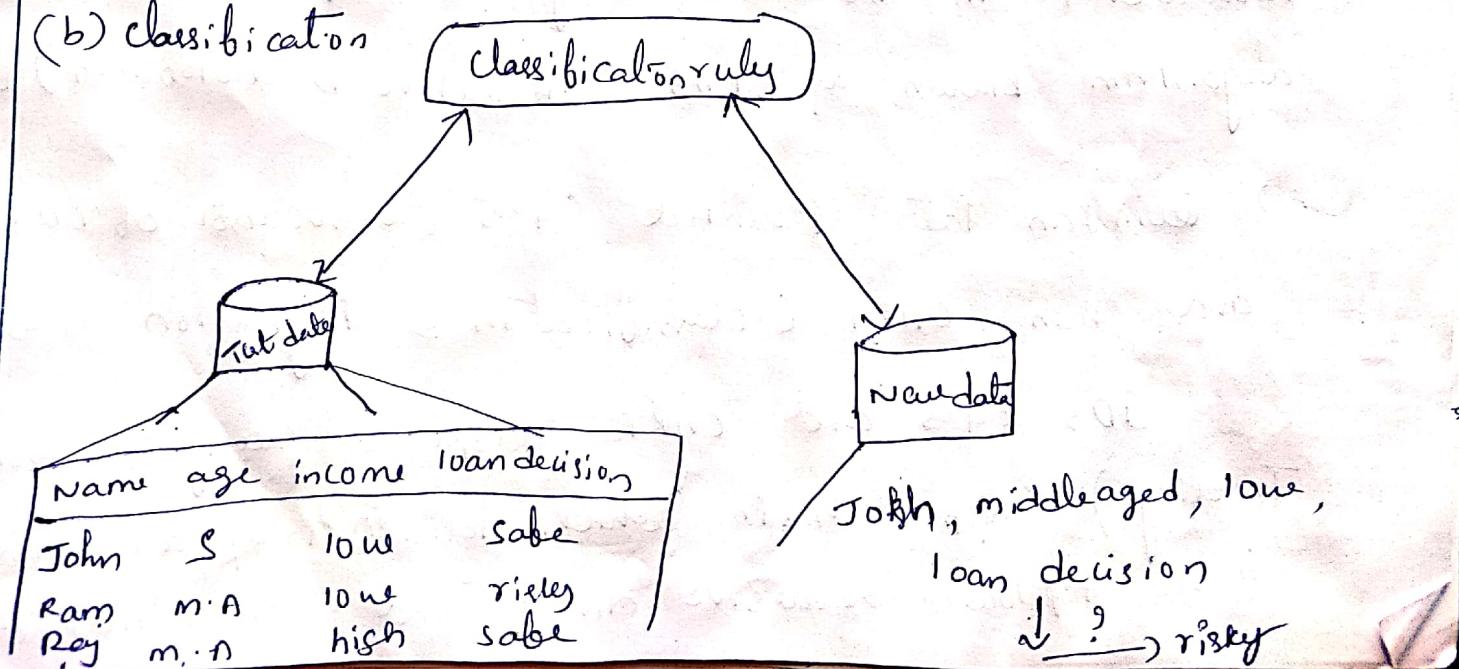
Data classification is two-step process

- ① Learning Step :- where a classification model is constructed. In first step, a classifier is built describing a pre-determined set of data classes (or) concepts. This is learning step (Training phase), where a classification algorithm builds classifier by analyzing (or) learning from a training set made up of database tuples & their class labels.
- ② Classification Step. Test data are used to estimate the accuracy of the classification rules. If accuracy is considered acceptable, the rules can be applied to the classification of new data tuples.



(a) Learning step

(b) Classification



141

The data classification process (a) learning : Training data are analyzed by a classification algorithm

(b) classification : Test data are used to estimate the accuracy of the classification rule.

→ Because the class label of each training tuple is provided this step is also known as supervised learning that is the learning of the classes is supervised in that it is told to which class each training tuple belongs

unsupervised learning is (clustering) in which class label of each training tuple is not known.

Decision Tree Induction

→ During late 1970's and early 1980's J. Ross Quinlan, researcher in machine learning, developed a decision tree algorithm known as ID3 (Iterative Dichotomiser).

→ Quinlan later presented C4.5 (a successor of ID3) and then CART (Classification & Regression Trees).

ID3, C4.5 and CART adopt a greedy (ie Non back tracking) approach in which decision trees are constructed in a topdown recursive divide and conquer manner.

Decision Tree Induction

(3)

Tree " Decision Tree induction is the learning of decision from class labeled training tuple .

" A Decision Tree is a flow-chart like tree structure where each internal node denotes a test on an attribute, branch represents an outcome of the test, and each leaf node holds a class label . The top most node is a tree is the root node .

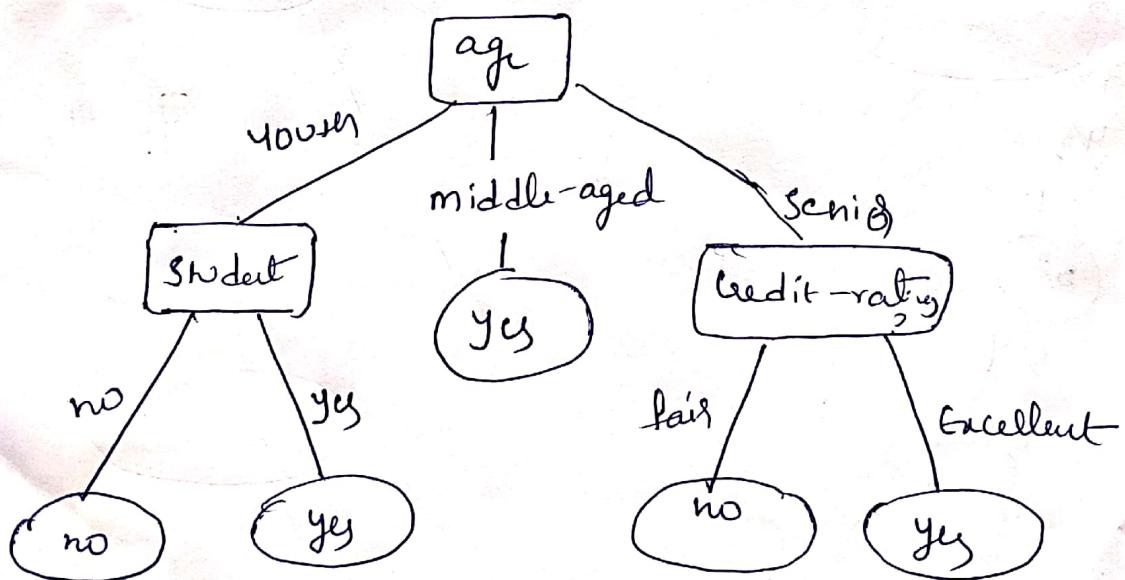


fig Decision Tree for the concept buys-computer

Selection measures

(4)

most decision Tree algorithms follow a top-down approach, which starts with a Training set of tuples and their associated class labels.

Algorithm: Generate decision - Tree. Generate a decision Tree from the training tuples of data partition, D.

Input:

- ✓ Data partition, D, which is set of training tuples and their associated class labels.
- ✓ Attribute list , the set of candidate attributes.
- ✓ Attribute - Selection method :- the splitting criterion that "best" partitions the data tuples into individual classes.

Output: A Decision Tree.

method:

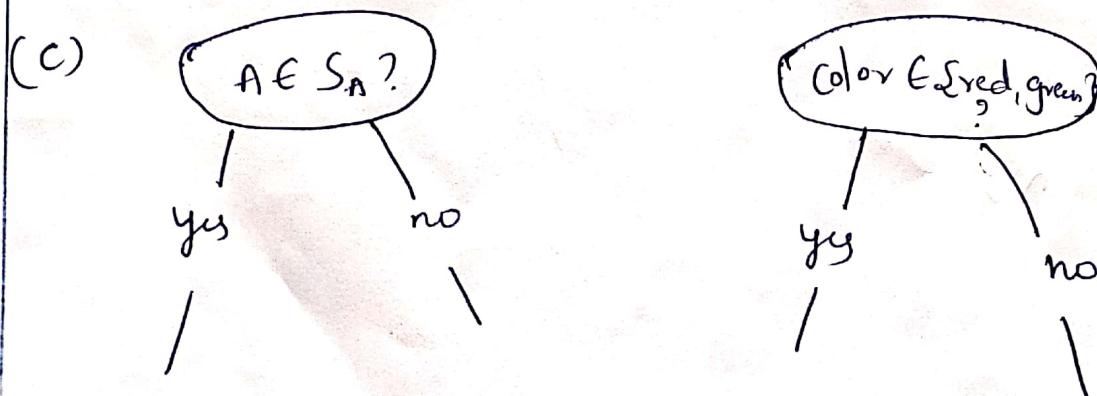
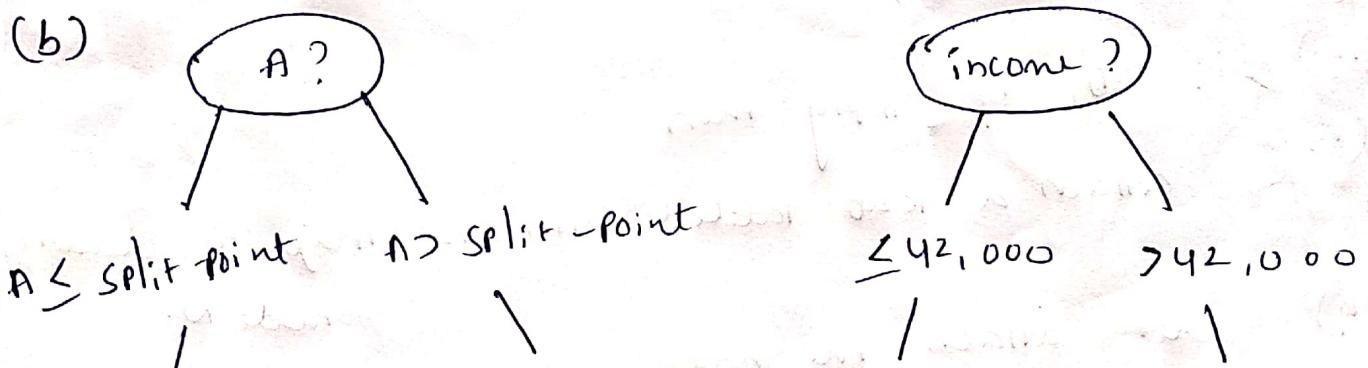
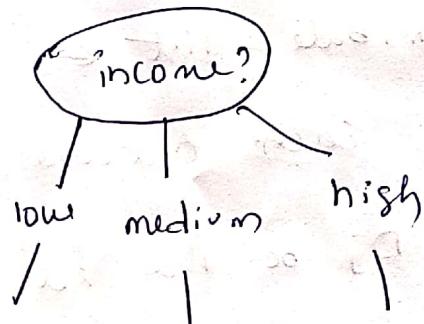
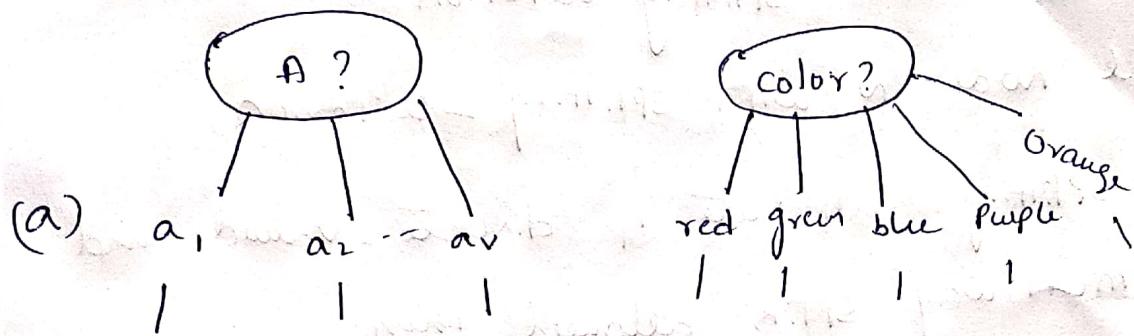
- (1) Create a node N;
- (2) if tuples in D are all of the same class c, then return N as a leaf node labeled with the class C;
- (3) if attribute-list is empty then return N as a leaf node labeled with the majority class in D;

(5)

- (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
- (9) attribute-list \leftarrow attribute-list - splitting attribute;
- (10) for each outcome j of splitting-criterion
- (11) let D_j be the set of datatuples in D satisfying outcome j ;
- (12) if D_j is empty then
- (13) attach a leaf labeled with the majority class in D to node n .
- (14) else attach the node returned by Generate-Decision-Tree to node n ;
 end for
- (15) return n ;

Splitting criterion

Example



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. Conceptually, the "best" splitting criterion is the one that most closely results in such a scenario.
- Attribute selection measures are also known as splitting rules because they determine how the tuples at a given node all have to be split.
- The attribute selection measure provides a ranking for each attribute describing the given training tuple.
- The attribute having the best score for the measure is chosen as the splitting attribute for the given tuple.
- Attribute selection measures describe 3 types
 - ① Information gain
 - ② Gain ratio
 - ③ Gini index

Q1

Information gain :- ID3 uses information gain as its attribute

Selection measure

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

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

The term $\frac{|\mathcal{D}_j|}{|\mathcal{D}|}$ acts as the weight of the jth partition.

$Info_A(D)$ is the expected information required to classify a tuple from D based on the partitioning by A.

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

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

Example

RID	age	income	student	credit-rating	class: buys-computer
1	young	high	no	fair	no
2	y	h	no	excellent	yes
3	middle-aged	h	no	fair	yes
4	senior	medium	no	fair	yes
5	s	low	yes	fair	no
6	s	l	yes	excellent	yes
7	m	l	yes	excellent	yes
8	y	l	yes	excellent	yes
9	y	medium	no	poor	no
10	s	low	yes	poor	yes
11	y	medium	yes	poor	yes
12	m	m	yes	poor	yes
13	m	high	no	poor	yes
14	s	medium	yes	poor	yes
			no	poor	no

Q1

distinct classes (i.e; m=2) & 4s, no³

2 (7)

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

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

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

$$\begin{aligned} \text{Gain}(\text{age}) &= \text{Info}(D) - \text{Info}_{\text{age}}(D) = 0.940 - 0.694 \\ &= 0.264 \text{ bits.} \end{aligned}$$

By $\text{Gain}(\text{income}) = 0.029 \text{ bits}$, $\text{Gain}(\text{credit-rat}) = 0.048 \text{ bits}$

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

Because age has the highest information gain among the attributes. it is selected as the splitting attribute.

node n is labeled with age and branches are grown for each of the attribute's values.

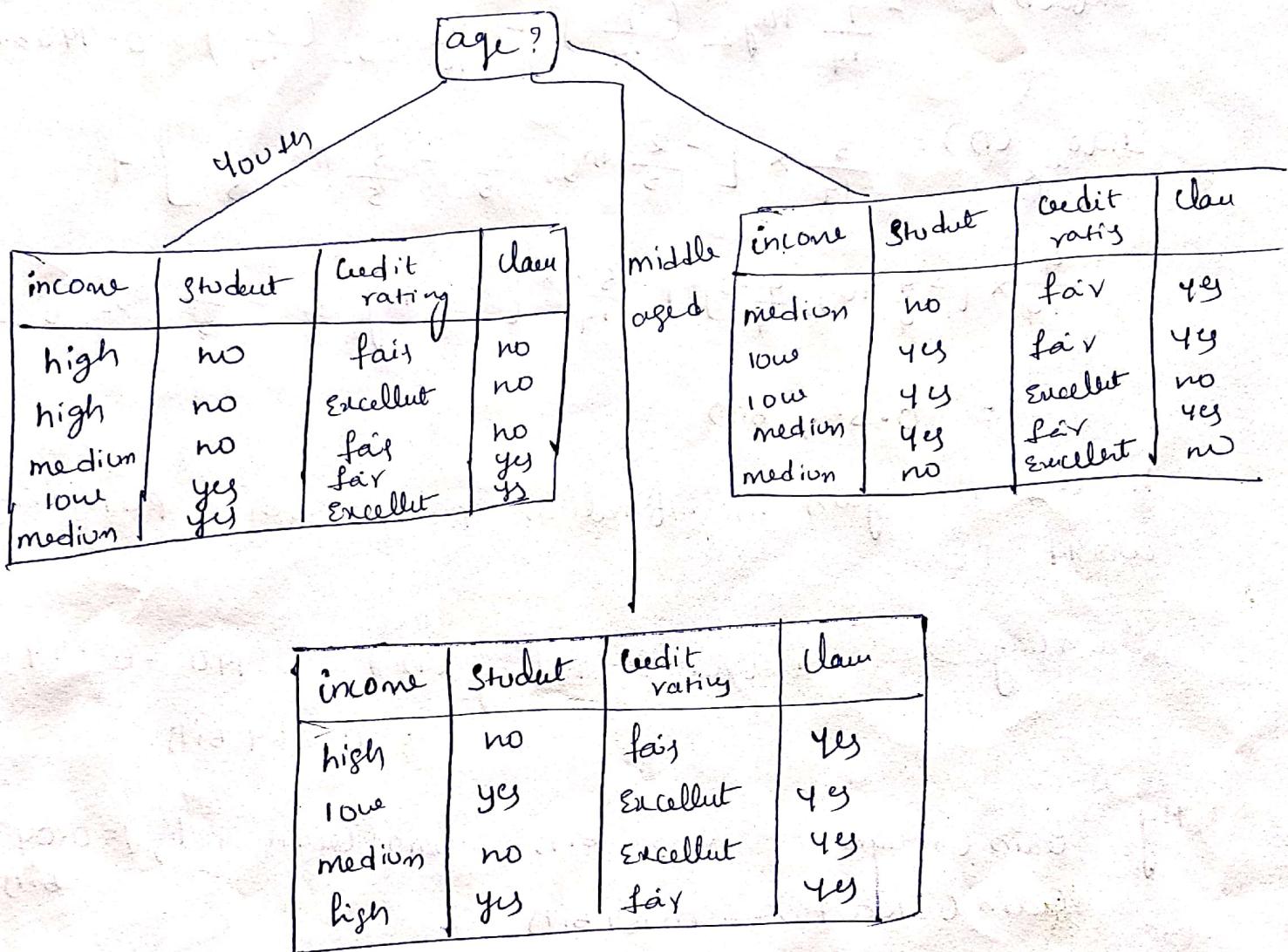


fig: decision tree.

Other Attribute Selection measure

- ① MDL (minimum description length) principle have the least bias toward multivalued attribute.
 - MDL-based measure use encoding techniques to define the "best" decision tree as the one that requires the fewest number of bits to both
 - (1) encode the Tree
 - (2) encode the exceptions to the tree
- ② multivariate splits : where the partitioning of tuples is based on a combination of attribute rather than a single attribute.
- CART system, for example, can find multivariate splits based on a linear combination of attribute
- multivariate splits are a form of attribute construction (feature) where new attribute are added to the existing ones.

Gain Ratio :- The information gain measure is biased for tests with many outcomes.

C4.5, a successor of ID3, uses an extension to information gain known as gain ratio, which attempts to overcome this bias. It applies a kind of normalization to information gain using a "split information" value defined analogously with Info(D) as

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

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

computation of gain ratio for the attribute income

A test on income splits the data of Table (Previous to) into three partitions, namely low, medium & high, containing four, six and four tuples. To compute the gain ratio of income

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

$$\text{we have } \text{Gain}(\text{income}) = 1.557 \\ 0.029 \therefore \text{GainRatio}(\text{income}) = \frac{0.029}{1.557} = 0.019$$

(9)

Gini Index

The Gini index is used in CART. using the notation described above. the Gini index measures the impurity of D , a data partition (δ) set of Training tuples as

$$Gini(D) = 1 - \sum_{i=1}^m p_i^2 \quad (\text{It is measure of the impurity (inequality) of } D)$$

where p_i is the probability that a tuple in D belongs to class C_i and is estimated by $|C_i \cap D| / |D|$

The Gini index considers a binary split for each attribute. Let's first consider the case where A is a discrete-valued attribute having v distinct values $\{a_1, a_2, \dots, a_v\}$, occurring in D . To determine binary split on A , we examine all of possible subsets that test for attribute A of the form " $A \in S_A?$ "

Ex: if income has three possible values $\{\text{low}, \text{medium}, \text{high}\}$ then possible subsets are $\{\text{low}, \text{medium}, \text{high}\}$, $\{\text{low}, \text{medium}\}$, $\{\text{low}, \text{high}\}$, $\{\text{medium}, \text{high}\}$, $\{\text{low}\}$, $\{\text{medium}\}$, $\{\text{high}\}$ and $\{\}$

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

If a binary split on income, partitions D into D_1 & D_2 the gini index of D given that partitioning is

(19)

For continuous-valued attribute

$$\Delta \text{Gini}(A) = \text{Gini}(D) - \text{Gini}_{\text{part}}(D)$$

Example Let D be the Training data, where there are 9 tuples belonging to the class buys-computer = yes and the remaining five tuples belong to buys-computer = no.

Gini index

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

Consider subset $\{\text{low, medium}\}$. This would result in 10 tuples in partition D_1 , satisfying the condition "income $\in [\text{low, medium}]$ ". The remaining four tuples of D would be assigned to partition D_2 .

$$\text{Gini}_{\text{income} \in \{\text{low, medium}\}}(D)$$

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

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

$$= 0.450$$

$$= \text{Gini}_{\text{income} \in \{\text{high}\}}(D)$$

My Gini index values for splits on the remaining subsets are 0.315 for subsets $\{\text{low, high}\}$ & $\{\text{medium}\}$ and 0.300 for

the subsets {medium, high} and low. Therefore the best binary split for attribute income is on {medium, high} & {low} because it minimizes Gini index.

Evaluating the attribute, we obtain {youth, senior} & {middle-aged} as the best split for age with Gini index of 0.375; the attributes {student} and {credit-rating} are both binary, with Gini index values of 0.367 & 0.429.

Find reduction in impurity using...

$$\text{income } \text{Gini}(A) = \text{Gini}(D) - \text{Gini}_A(D)$$

$$= 0.459 - 0.300 = 0.159$$

Step1: compute the impurity of D

$$\text{Gini(age)} = 0.459 - 0.315 = 0.144$$

$$\text{Gini}_{\text{student}} \quad \text{Binary} = 0.459 - 0.367 = 0.092$$

$$\text{Gini}_{\text{credit-rating}} \quad \text{Binary} = 0.459 - 0.429 = 0.03$$

Income is selected with minimum gini index & highest reduction in impurity.

Four Steps:- of Gini index

(A)

- ① Find the impurity of D, using formula i
- ② find the impurity of each resulting partition using

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

- ③ find reduction in impurity using formula

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

- ④ now select the best attribute which gives the minimum gini index

Bayes classification method

(11)

- Bayesian classification (8) Bayesian classifiers are statistical classifiers. They can predict class member probabilities such as the probability that a given tuple belongs to a particular class.
- Bayesian classification is based on Baye's theorem, Bayesian classifier is also known as Naive Bayesian classifier.
- In Bayesian setting, x is considered "evidence". As usual, it is described by measurements made on a set of n attributes.
- Let H be some hypothesis such as that the data tuple x belongs to a specified class c . For classification problem we want determine $P(H|x)$
 $P(H|x)$ is the posterior probability. Suppose our world of data tuple is confined to customers described by the attributes age & income respectively and that x is 35-year old customer with an income of \$40,000. H is the hypothesis that our customer will buy a 'computer' given age & income.

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

My $P(X|H)$ is the posterior probability of X conditioned on H . i.e. Probability that a customer x , is 35 years old and earning 40,000 given we know the customer will buy a computer.

Bayesian classification

(12)

Bayesian classification is based on Baye's theorem.

Bayesian classification is also known as naive Bayesian classification.

Bayes Theorem:

$P(H|z)$ is the posterior probability of a posterior probability, of H conditioned on x .

$$P(H|z) = \frac{P(z|H) P(H)}{P(z)}$$

$$P(c_i|z) = \frac{P(z|c_i) P(c_i)}{P(z)}$$

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

$x = \text{Age} = \text{youth}, \text{income} = \text{medium}, \text{Student} = \text{yes}$,

$\text{Credit rating} = \text{fair}$)

$P(z|c_i) P(c_i)$ for $i = 1, 2, \dots, m$, $P(c_i)$ the prior probability of each class.

$$P(\text{buys - comp} = \text{yes}) = 9/14 = 0.643$$

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

To compute $P(x|c_i)$ for $i=1, 2$

$$P(\text{age} = \text{youth} | \text{buys - comp} = \text{yes}) = 2/9 = 0.222$$

$$P(\text{age} = \text{youth} | \text{buys - comp} = \text{no}) = 3/5 = 0.600$$

$$P(\text{income} = \text{medium} | \text{buys - comp} = \text{yes}) = 4/9 = 0.444$$

$$P(\text{income} = \text{medium} | \text{buys - comp} = \text{no}) = 2/5 = 0.400$$

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

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

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

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

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

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

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

$$= 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$$

16y

(13)

$$P(x | \text{buys-comp} = \text{no}) = 0.660 \times 0.400 \times 0.200 \times 0.400 = 0.019.$$

$$P(x | \text{buys-comp} = \text{yes}) P(\text{buys-comp} = \text{yes}) = 0.044 \times 0.643$$

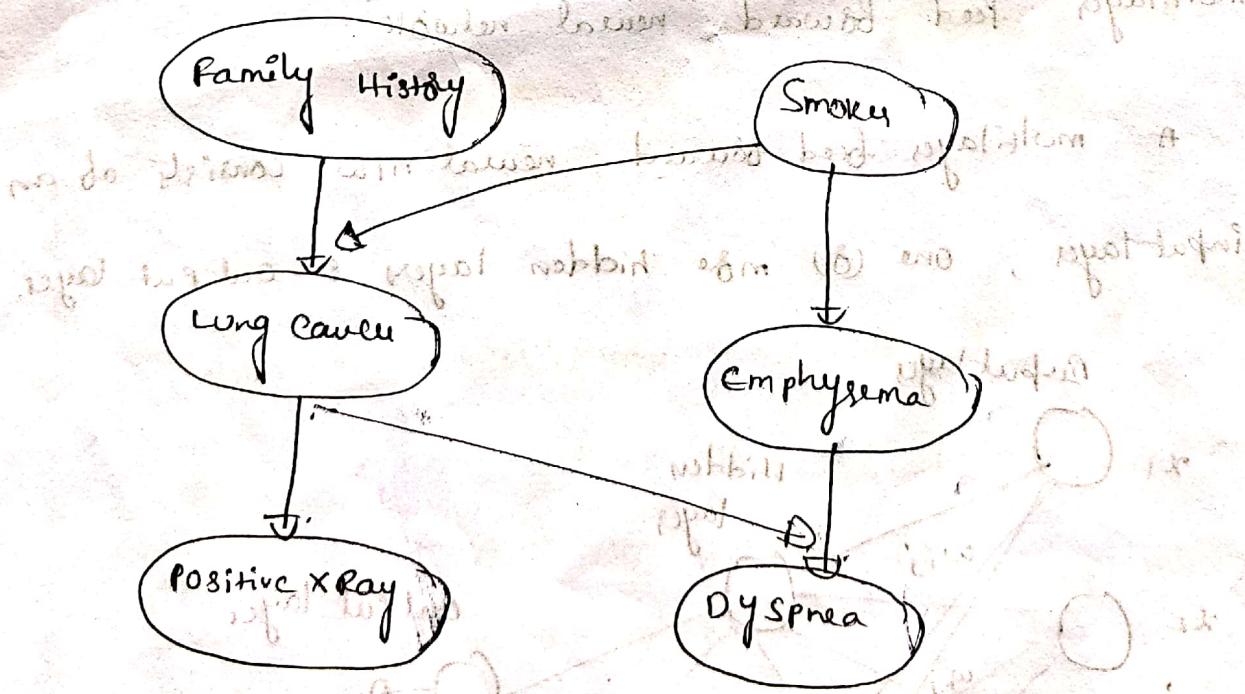
$$= 0.028$$

$$P(x | \text{buys-comp} = \text{no}) P(\text{buys-comp} = \text{no}) = 0.019 \times 0.357$$

$$= 0.007$$

naive bayesian classifier predicts $\text{buys-comp} = \text{yes}$ for tuple x .

A simple Bayesian belief network



Classification by Back propagation

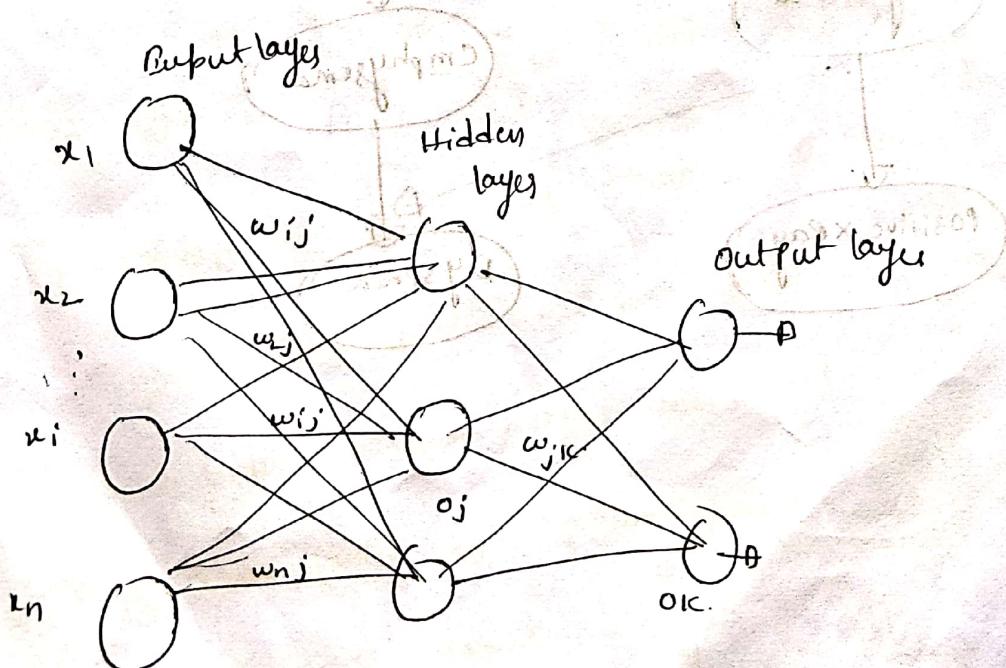
Back propagation is a neural network learning algorithm.

"Neural network is a set of connected input/output units in which each connection has 'weight' associated with it".

* A multilayer feed-forward neural network

* The back propagation algorithm performing learning on a multilayer feed-forward neural network.

* A multilayer feed forward neural nw consists of an input layer, one or more hidden layers & output layer.



Tree Pruning

- when a decision tree is built, many of the branches will reflect anomalies in the training data due to noise (or) outliers.
- Tree Pruning method addresses this problem of overfitting the data.
- unpruned tree tend to larger & complex but pruned tree tend to smaller and less complex easier to comprehend.

There are two common approaches to Tree Pruning

- ① Pre pruning
- ② Post pruning

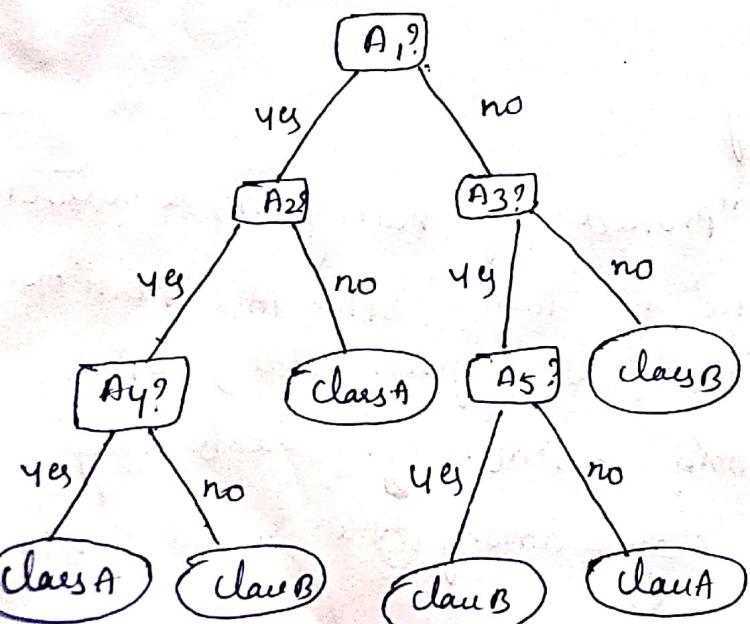
- Prepruning :- → A tree is "pruned" by halting its construction early (ex: by deciding not to further split (or) partition the subset of training tuples at a given node) upon halting, the node become a leaf. The leaf may hold the most frequent class among the subset tuples.
- when constructing a tree measure such as Information gain, Gini index etc. can be used to assess the goodness of split.

Ib Partitioning the tuple at a node would result in a split that falls below a prespecified threshold; then further partitioning of the given subset is halted.

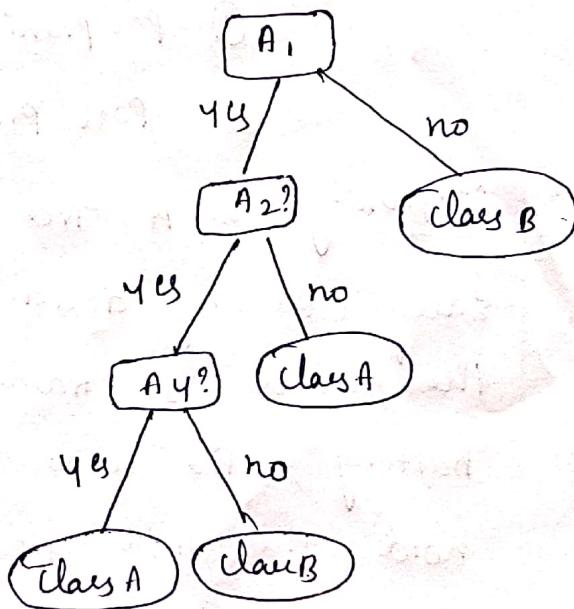
→ High thresholds could result in oversimplified trees, while low thresholds could result in very little simplification.

Post pruning, which removes subtrees from a "fully grown" tree. A subtree at a given node is pruned by removing its branches & replacing with its leaf.

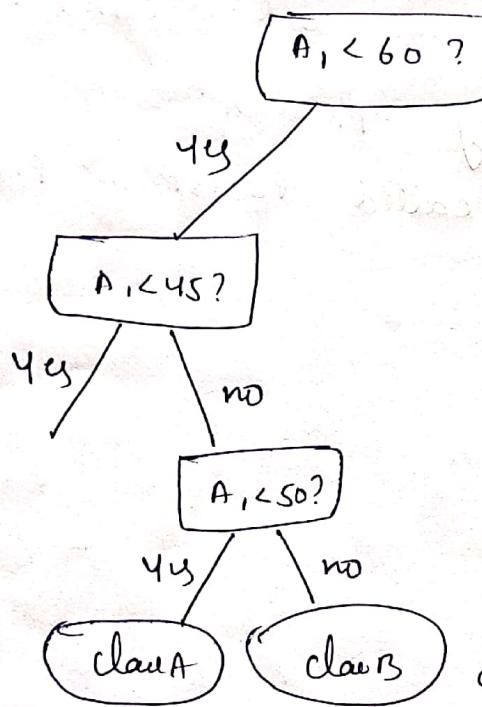
→ The leaf is labeled with the most frequent class among the subtree being replaced. For example A_3 in the unpruned tree of below figure.



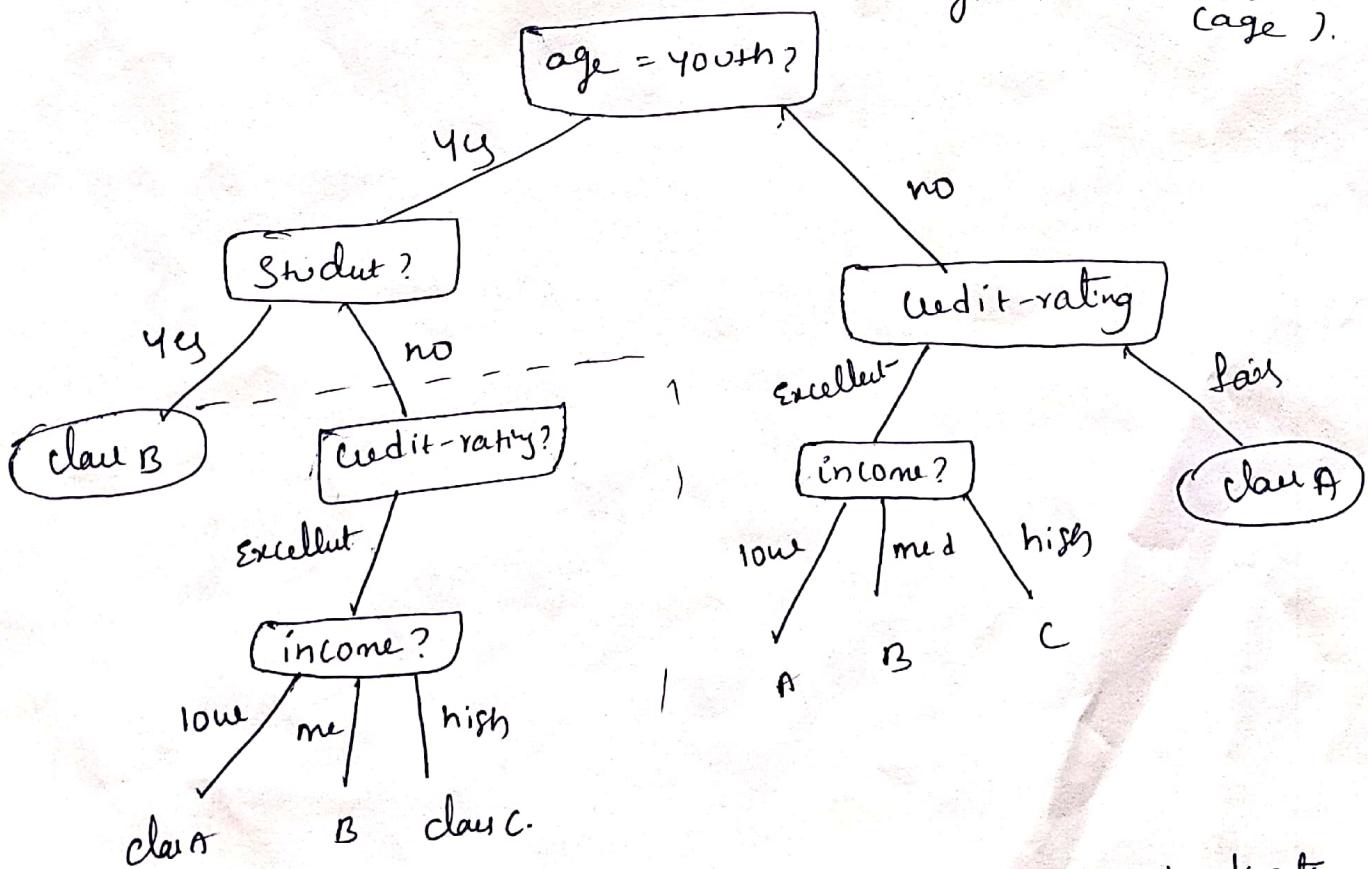
fig(a) unpruned decision tree



fig(b) Pruned decision tree



(a) Subtree repetition, where an attribute is repeatedly tested along a given branch of Tree (age).



(b) Subtree replication, where duplicate subtree exist within a tree (C.R.)

- Cost complexity pruning algorithm used in CART is an example of the post pruning approach.
- C4.5 uses a method called Pessimistic Pruning.
Alternatively

Scalability and Decision Tree Induction

(16)

"How scalable is decision tree induction"? The efficiency of existing decision tree algorithms, such as ID3, C4.5 & CART has been well established for small datasets.

- Efficiency becomes an issue of concern when these algorithms are applied to the mining of very large real world databases.
- In data mining applications, very large training set of millions of tuples are common. The training data will not fit in memory. Therefore decision tree construction becomes inefficient due to swapping of the training tuple in and out of main & cache memory.
- Several scalable decision tree induction methods have been introduced in recent studies. Rainforest, for example adapts to the amount of main memory available and applies to any decision tree induction algorithm.
- This method maintains a AVC-set (Attribute value class label) for each attribute

}

(17)

age	buys-computer	
	Yes	No
Youth	2	3
middle-aged	4	0
senior	3	2

income	buy-computer	
	Yes	No
low	3	1
medium	4	2
high	2	2

Student	buys-computer	
	Yes	No
Yes	6	1
No	3	4

credit-rating	buys-computer	
	Yes	No
fair	6	2
excellent	3	3

fig: The use of data structures to hold aggregate information regarding the Training data (These Arc-sets describing previous Training data set) are one approach to improving Scalability of decision tree induction.

- "Boat" (Bootstrapped optimistic Algorithm for Tree construction) is a decision tree algorithm that takes a completely different approach to scalability
- It is not based on the use of any special data structure. Instead, it uses a statistical technique known as "bootstrapping" to create several smaller samples of the given training data, each of which fits in memory

- BOAT can use any attribute selection measure that selects binary splits and that is based on the notion of purity of partitions such as the Gini index.
- BOAT was found to be two to three times faster than Random Forest. Advantage of BOAT is that it can be used for incremental update. That is BOAT can take new insertions and deletions for the training data and update the decision tree to reflect these changes, without having to reconstruct the tree from scratch.

Visual mining for Decision Tree Induction

- How can visualize the Decision Trees means one of the method like perception-based classification (PBC).
- PBC is an interactive approach based on multidimensional visualization technique and allows the user to incorporate background knowledge about the data when building a decision tree.
- The resulting trees tend to be smaller than those built using Traditional decision tree induction methods.

"How can the data be visualized to support interactive decision tree construction?

- PBC uses a pixel-oriented approach to view multidimensional data with its class label information.
- The circle segments approach is adapted, which maps d -dimensional data objects to a circle that is partitioned into d segments, each subtending one attribute. Here an attribute value of an object is mapped to one colored pixel of each data object.
- PBC system displays a split screen, consisting of Data Interaction window and a knowledge Interaction window
Data Interaction window: displays the circle segments of the data under examination.
knowledge Interaction window:- displays the decision tree constructed.

Initially the complete training set is visualized in the Data Interaction window, while the knowledge Interaction window displays an empty decision tree

(18)

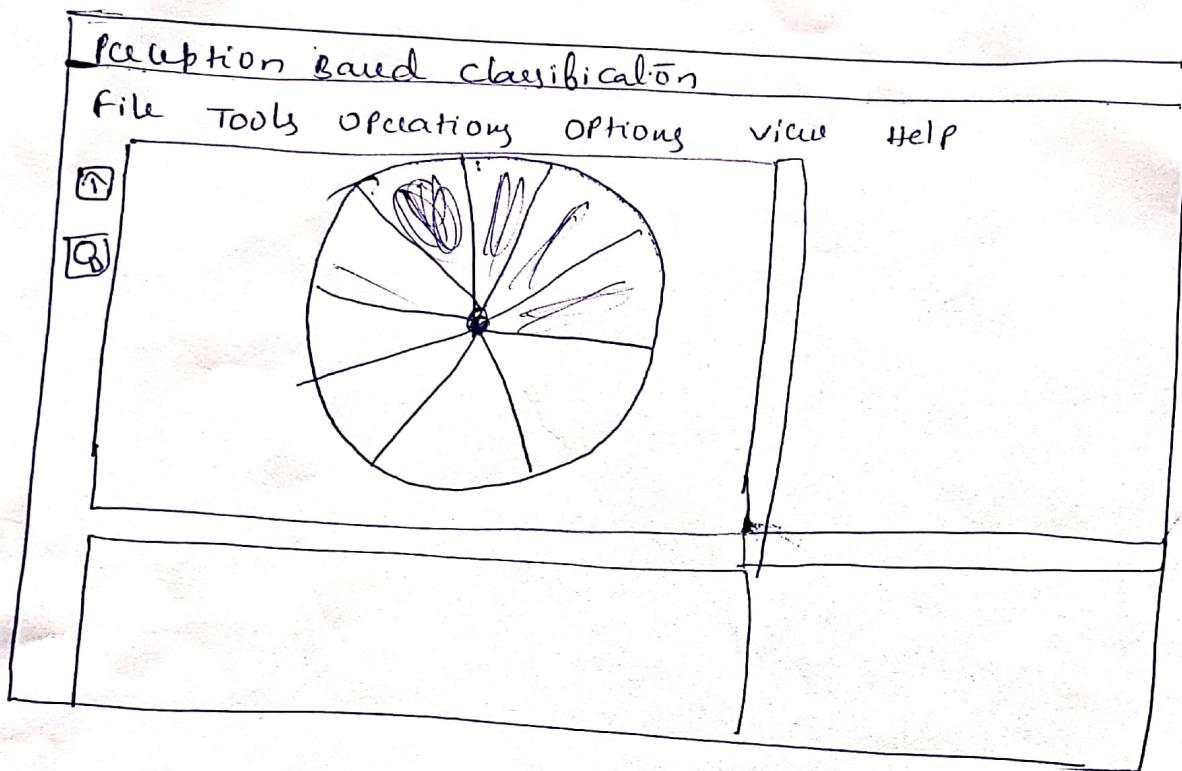


Fig : A screenshot of PBC.js interactive decision tree construction

Model evaluation and Selection

(19)

model selection: choosing one classifier over another

Metric for evaluating classifier performance

This section presents measures for assessing how good & how "accurate" your classifier is at predicting the class label of tuples.

measure

formula

Accuracy, Recognition rate

$$\frac{TP + TN}{P + N}$$

Error rate, misclassification rate

$$\frac{FP + FN}{P + N}$$

Sensitivity, True positive rate

$$\frac{TP}{P}; \frac{TN}{N}$$

Specificity.

Precision

$$\frac{TP}{TP + FP}$$

F, F₁, F-score

harmonic mean of precision &
recall

$$\frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

F_B, where B is a non-negative
real number

$$\frac{(1+B)^2 \times \text{Precision} \times \text{Recall}}{B^2 \times \text{Precision} + \text{Recall}}$$

There are the evaluation measures. TP , TN , FP , P , N refers to the number of true positive, true negative, false positive, positive, and negative samples.

There are four additional terms we need to know that are the "building blocks" used in computing many evaluation measures.

True Positive (TP):- There refers to the positive tuples that were correctly labeled by the classifier. Let TP be the number of true positives.

True negatives :- (TN) There are the negative tuples that were correctly labeled by the classifier. Let TN be the number of true negatives.

False Positive (FP): There are the negative tuples that were incorrectly labeled as positive. Let FP be the number of false positives.

False Negatives (FN): There are the positive tuples that were mislabeled as negative. Let FN be the number of false negatives.

These terms are summarized in the confusion matrix.

- based classification:-

20

Predicted class

Actual class

		Yes	No	Total
Yes	TP	FN	P	
	FP	TN		
Total	P'	N'	P + N	

Confusion matrix, with totals for positive & negative tuples.

Classes	Buy - com = Yes		Buy - com = No		Total	Recognition
	buys - comp = yes	buys - comp = no	buys - comp = yes	buys - comp = no		
buy - com = yes	6954	412	46	2588	7006	99.34
buy - com = no					3000	86.27
Total	7366		2634		10,000	95.42

Fig: Confusion matrix for the class buys - computer = Yes

2) buys computer = NO .

Example: Sensitivity and Specificity:- confusion matrix for medical data where the class values are Yes and no for a class label attribute - cancer.

classes	yes	no	Total	Recognition (%)
yes	T_p 90	F_n 210	300	30.00
no	F_p 140	T_n 9560	9700	98.56

The sensitivity of the classifier is $\frac{90}{300} = 30.00\%$.

The specificity is $\frac{9560}{9700} = 98.56$

The classifier overall accuracy = $\frac{9650}{10,000} = 96.50\%$.

The precision & recall measures are also widely used in classification. precision can be thought of as a measure of exactness (ie what percentage of truly labeled as positive are actually), recall is a measure of completeness .

$$\text{Precision} = \frac{T_p}{T_p + F_p}$$

$$\text{Recall} = \frac{T_p}{T_p + F_n} = \frac{T_p}{P}$$

Rule - Based classification:-

Example

$$\frac{TP}{TP+FP} = \frac{90}{90+140} = 39.13\%$$

$$\text{Precision} = \frac{90}{230} = 39.13\%$$

$$\text{Recall} = \frac{90}{300} = 30.00\%$$

In addition to accuracy-based measures, classification can also compare with respect to the following additional aspects.

Sensitivity:

Example

		Cancer = Yes	Cancer = No	Total
Actual class	Predicted class			
		Cancer = Yes	Cancer = No	Total
Cancer = Yes	Cancer = Yes	90	210	300
Cancer = No	Cancer = No	140	9560	9700
Total		230	9770	10000

$$\text{Sensitivity} = 30.0$$

$$\text{Recall} = 30.00\%$$

$$\text{Specificity} = 98.56$$

$$\text{Accuracy} = 96.40$$

$$\text{Precision} = 39.13$$

based on above

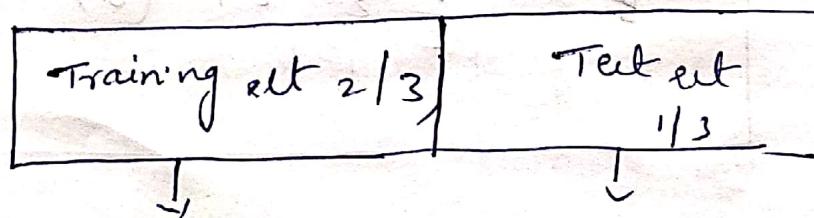
formulas.

Methods to evaluating classifier accuracy

- ① Hold out
- ② Random Sampling
- ③ Cross validation
- ④ Bootstrap

Hold out :-> In this method, the given data are randomly partitioned into two independent sets, like a Training dataset and test set.

- Two-thirds ($2/3$) of data are collected to the Training set, and the remaining one-third ($1/3$) is allocated to the test set.
- The Training data set is used to derive the model.
- The model's accuracy is estimated with the test set.



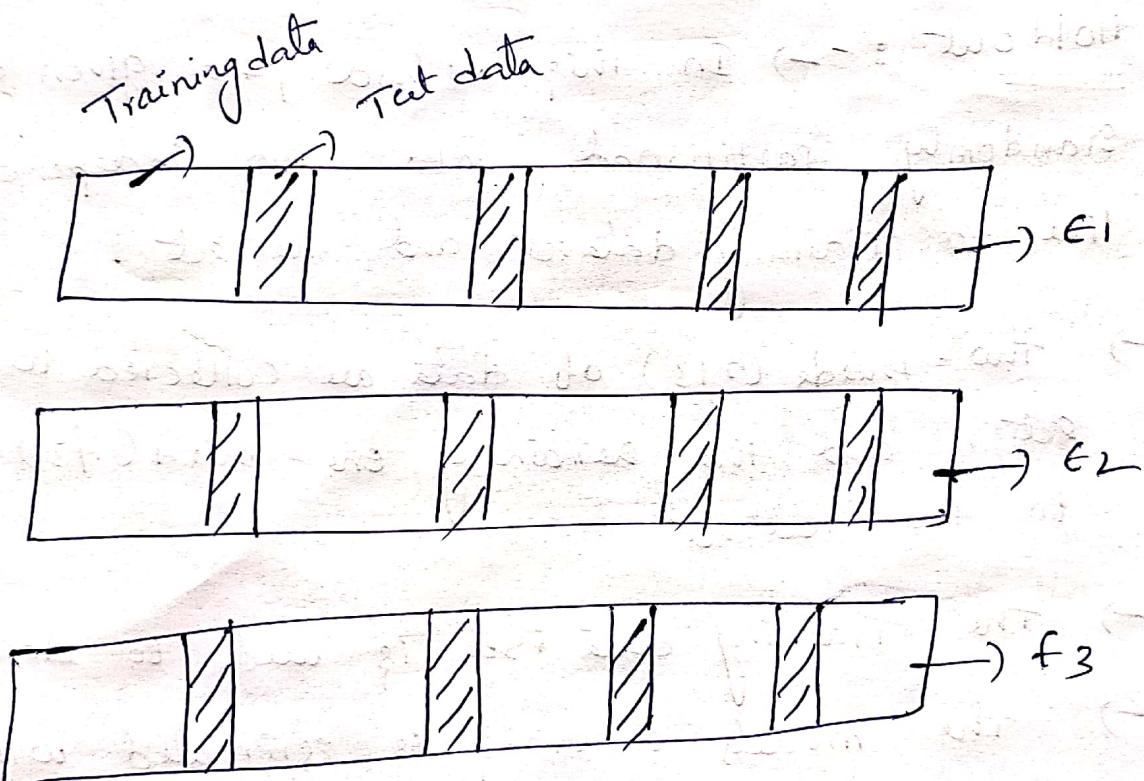
construct model

for accuracy estimations

Fig: Hold out method.

Random Sampling :- (2) Repeated Hold out method

- Random Sampling is a variation of the holdout method in which the holdout method is repeated K times.
- The overall accuracy estimate is taken as the average of the accuracies obtained from each iteration.



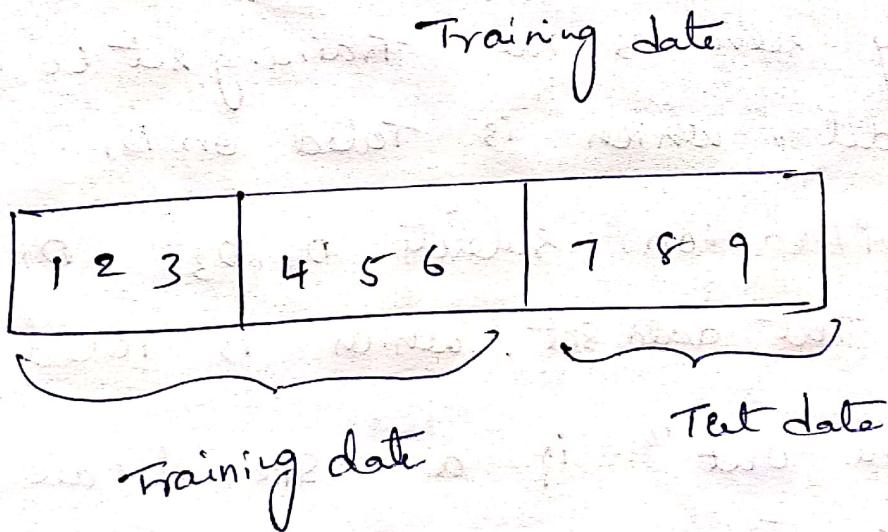
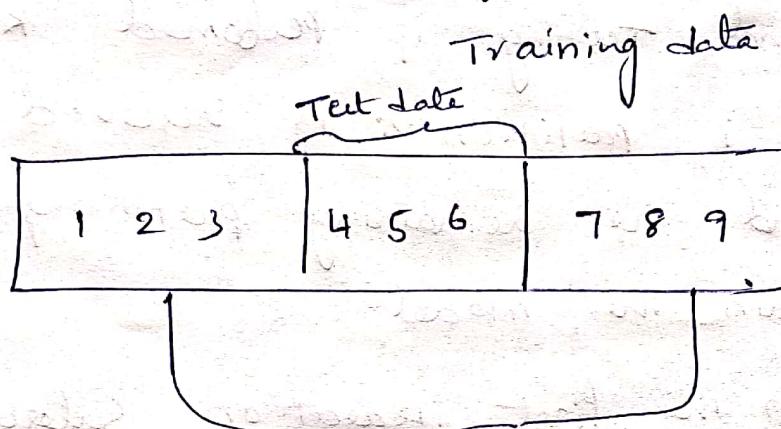
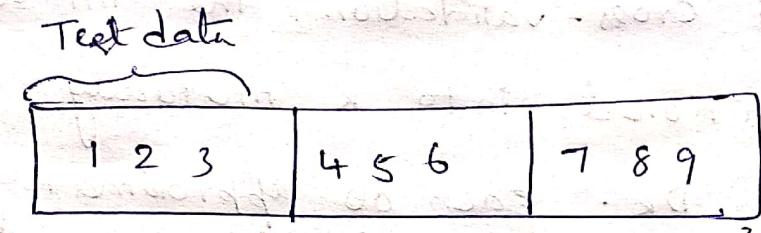
$$E = \frac{1}{K} \sum_{i=1}^K E_i \quad (2) \quad \frac{E_1 + E_2 + E_3}{3}$$

Cross validation

23

- In k -fold cross-validation, the initial data are randomly partitioned into k mutually exclusive (\emptyset) folds $D_1, D_2 \dots D_k$. Each of approximately same size.
- Training and Testing is performed k times
In iteration i , partition D_i is reserved as the test set, and the remaining partitions are collectively used to train the model.
- That is, in the first iteration, subsets $D_2, D_3 \dots D_k$ collectively serve as the "Training set" to obtain a first model, which is Tested on D_1 .
- Second iteration, subsets $D_1, D_3 \dots D_k$ collectively serve as Test data set, which is Tested on D_2
- "Leave one-out" is a special case of k -fold cross validation where k is set to the number of initial tuples. That is only one sample is "left out" at a time for the test set.
- "Stratified cross-validation", the folders are stratified so, that the class distribution of tuples in each fold is approximately the same as in the initial data.

3-fold cross validation



Boot strap:-

24

- works well with small datasets.
- samples the given training tuple uniformly with replacement.
- each time tuple is selected, it is equally likely to be selected again and re-added to the Training set.
- Several Boot strap methods like .632 bootstrap.
A dataset with d tuples is sampled d times, with replacement, resulting in Training set of d samples. The data tuple that did not make it in to the Training set end up forming the test set.
- About 63.2% of the original data end up with bootstrap, and Remaining 36.8% from the test set
- $$(1 - 1/d)^d = e^{-1} \approx 0.368.$$
- Repeat this process k times, overall accuracy of the model.

$$ACC(m) = \frac{1}{k} \sum_{i=1}^k (0.632 \times ACC(m_i)_{\text{test}} + 0.368 \times ACC(m_i)_{\text{Train}})$$

Rule-Based classification

25 ①

201

- x Rule-Based classifier (81) classification is supported by a set of IF-THEN rules.
- x we will see how such rules are used for classification

Using If-Then rules are one of the rule-based classification method (81) learning method.

Syntax:

If condition
↓
Rule antecedent

Then conclusion
↓
Rule consequent

In rule (81) of if and then

Pre-condition

Ex:-

R₁: If age = youth AND income = high AND student = yes

then buys comp = yes

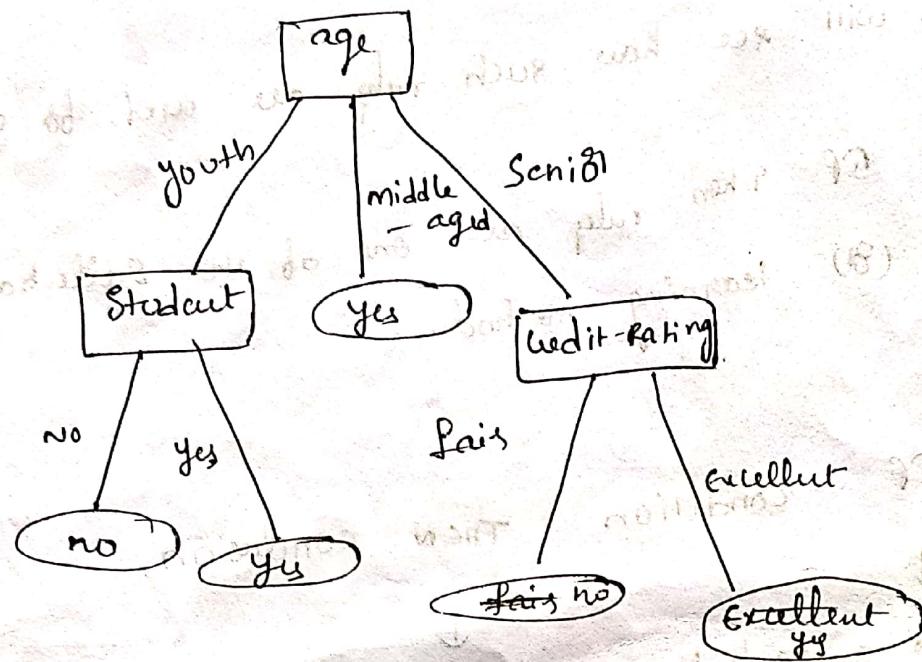
Another way rep IF-THEN Rule.

(age = youth) \wedge (income = high) \wedge (student = yes) \Rightarrow
buys comp = yes.

$$\text{Coverage}(R) = \frac{n(\text{Cover})}{|D|}, \quad \text{Accuracy}(R) = \frac{n(\text{Correct})}{n(\text{Cover})} \\ = \frac{214}{281} = 14.28\% \quad = \frac{212}{281} = 75.71\%$$

How to extract the Rule means: i.e. Rule are build (8) is extracted from the Decision tree. Decision tree is popular method to extract the rules.

Ex:



Here condition is taken from the path. Path is taken from root node & conclusion is taken from the leaf node.

R₁: age = youth AND Student = no then buys-comp = no

R₂: age = youth AND Student = yes then buys-comp = yes.

R₃: age = middle-aged then buys-comp = yes.

R₄: age = Senior AND CreditRating = fair then buyscomp = no

R₅: age = Senior AND CreditRating = excellent then
buys-comp = yes.

Age income student creditrating buys comp

y	h	NO	fair	yes
m	h	NO	excellent	NO
s	medium	NO	fair	yes
g	medium	NO	"	NO
y	low	yes	"	yes
T		NO		NO

If (age = youth) then

buycomp = yes

If age = youth

then buycomp = yes

If age = middle-aged

then buycomp = yes

**IF age = youth
THEN buycomp = yes**

**IF age = income = high
THEN buycomp = yes**

**IF income = high
AND creditrating = fair
THEN buycomp = yes**

If age = middle-aged
then buycomp = yes

Rule Induction using a Sequential covering Algorithms

- * If-THEN rules can be extracted directly from the training data (without having to generate a decision tree) using a sequential covering algo.
- * Sequential covering algos are the most widely used approach to mine classification rules.

algo:-

Input:

- ✓ D, a data set class-labeled tuples;
- ✓ Attr_val, the set of all attribute & their possible values.

Output:

A set of IF-THEN rules.

- * In this algo first rule learns sequentially i.e. it learns first rule next learns second rule. until all rules are completed.

method

- (1) Rule_Set = {}
- (2) for each class c do
- (3) repeat
- (4) Rule = Learn-one-Rule(D, Attr_val, c);
- (5) Remove tuples covered by Rule from D;
- (6) Rule_Set = Rule_Set + Rule;
- (7) until terminating condition;
- (8) return Rule_Set;

How are rules learned?

Rules are grown in a general-to-specific manner. first start with empty set. and then gradually keep appending attribute tests to it. suppose our training data set, D consists of loan application data. Attributes regarding each applicant include age, income, loan term, residence, credit rating. The classifying attribute is loan-decision, which indicates that whether a loan is accepted (A) or Rejected.

To learn a rule for the class "accept"

If THEN loan-decision = accept.

age	loan term	income	credit rating	loan decision
young	short	high	Excellent	Accept
middle-age	long	high	Excellent	Accept
young	medium	high	Fair	Accept
young	short	low	poor	reject

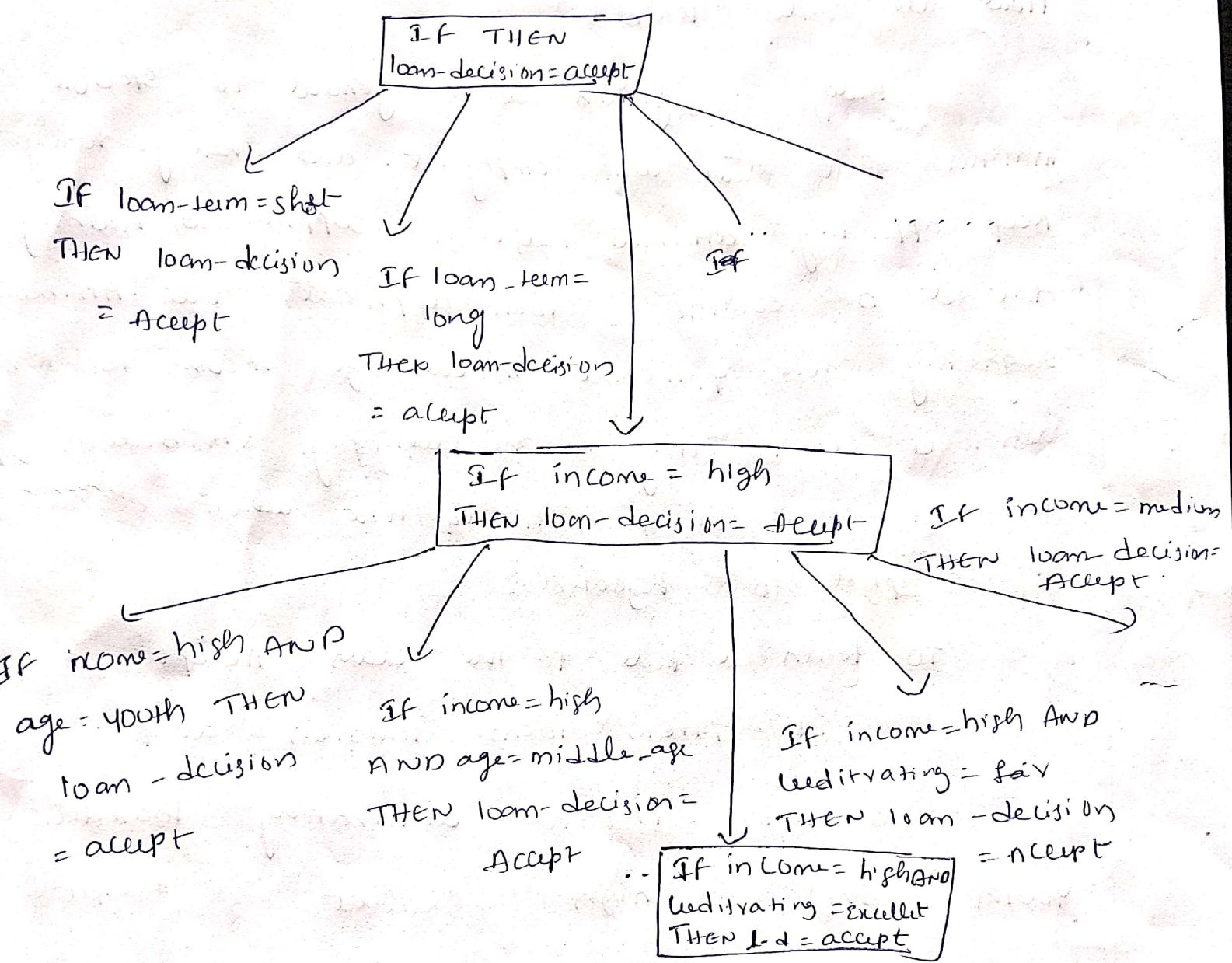


Fig.: A general-to-specific search through rule space.

Training data will contain many attributes, each of which may have several possible values. Finding an optimal rule set.

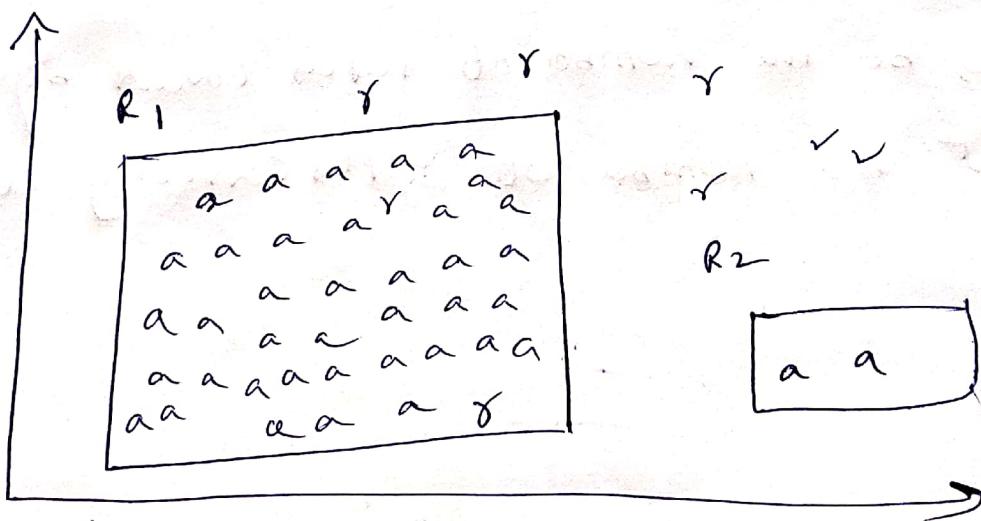
Rule Quality measures

Learn-one-Rule needs a measure of rule quality. Every time it considers an attribute test, it must check to see if appending such a test to the current values condition will result in an improved rule.

Choosing two rules based on Accuracy

Consider the two rules as illustrated below. Both are for the class loan decision = accept. we use 'a' to rep the tuples of class "accept", 'g' for the tuples of class reject.

Rule R_1 correctly classifies 38 of the 40 tuples it covers. Rule R_2 covers only two tuples, which it correctly classifies. Their accuracies are 95% & 100%, R_2 greater accuracy than R_1 .



Rule for the class loan decision = accept

Another measure is FOIL (First order inductive learner)

$$\text{FOIL-Gain} = \text{pos}' \times \left(\log_2 \frac{\text{pos}'}{\text{pos}' + \text{neg}'} - \log_2 \frac{\text{pos}}{\text{pos} + \text{neg}} \right)$$

Rule pruning

$$\text{FOIL-Prune}(R) = \frac{\text{pos}-\text{neg}}{\text{pos}+\text{neg}'}$$

- coverage is the percentage of tuples that are covered by the rule.
- accuracy - what percentage of them the rule can correctly classify.
- n_{covers} be the number of tuples covered by R
 n_{correct} the number of tuples correctly classified by R .

Introducing Ensemble methods

- An ensemble combines a series of K learned models m_1, m_2, \dots, m_K with the aim of creating an improved composite classification model. m^*
- A given dataset D , is used to create K training sets D_1, D_2, \dots, D_K is used to generate classifiers m_i .
- Given a new data tuple to classify, the base classifiers each vote by returning a class prediction.
- The ensemble returns a class prediction based on the votes of the base classifier.

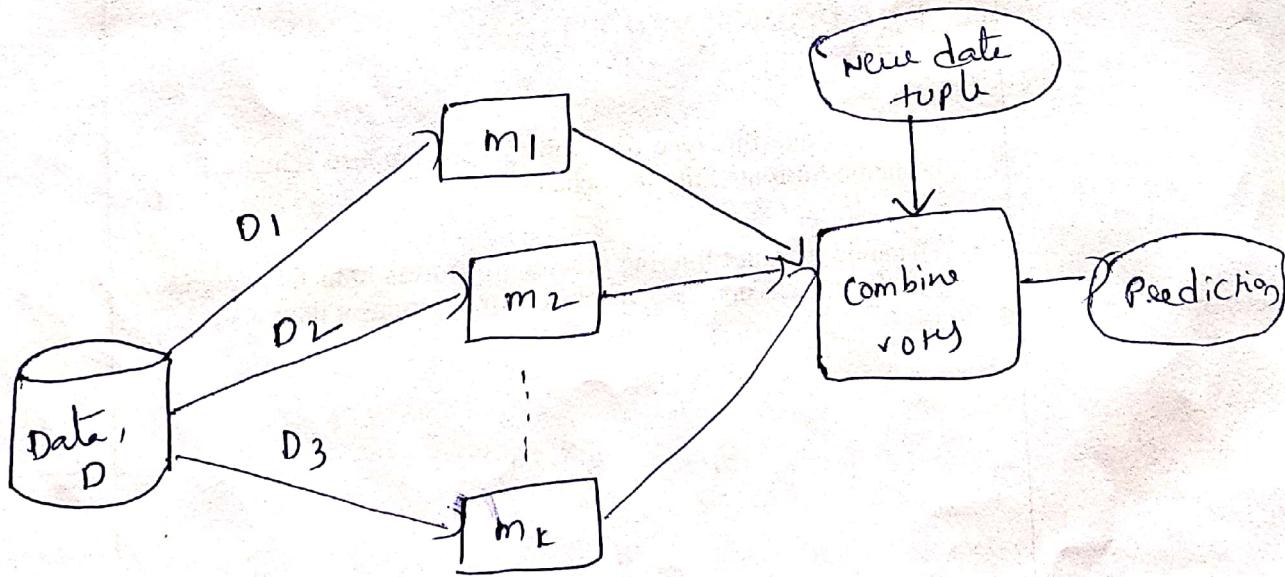


fig. Increasing classifier accuracy.

Ensemble methods are

30

- ① Bagging
- ② Boosting & Ada Boosting.
- ③ Random forests

- Bagging: → Bagging works as a method of increasing accuracy. Suppose that you are patient and would like to have diagnosis made based on your symptoms.
- Instead of asking one doctor, you may choose to ask several. Final diagnosis is made based on majority vote, where each doctor gets an equal vote.
 - now replace doctor by a classifier. combine a series of models to generate composite classification model.
 - Bagging is also called as bootstrap aggregation

Algorithm: The bagging algorithm - Create an ensemble of classification models for a learning scheme where each model gives an equally weighted prediction

- Input:
- D, set of d Training tuples
 - K, number of models in the ensemble
 - a classification learning scheme (Decision Tree algo, Naive Bayes algo)

Output: A composite model $m(x)$

method

- (1) for $i = 1$ to K do // make K models
- (2) make bootstrap sample D_i , by sampling D with replacement;
- (3) use D_i and the learning scheme to derive a model m_i ;
- (4) end for.

Boosting and Ada Boost

→ Suppose that as a patient, you have certain symptoms. Instead of consulting one doctor, you choose to consult several.

- Suppose you assign weights to the value (0, 1) worth of each doctor's diagnosis, based on the accuracies of previous diagnoses they have made.
- The final diagnosis is then a combination of the weighted diagnoses.
- In Boosting, weights are also assigned to each training tuple. A series of K classifiers is learned iteratively.

Bagging

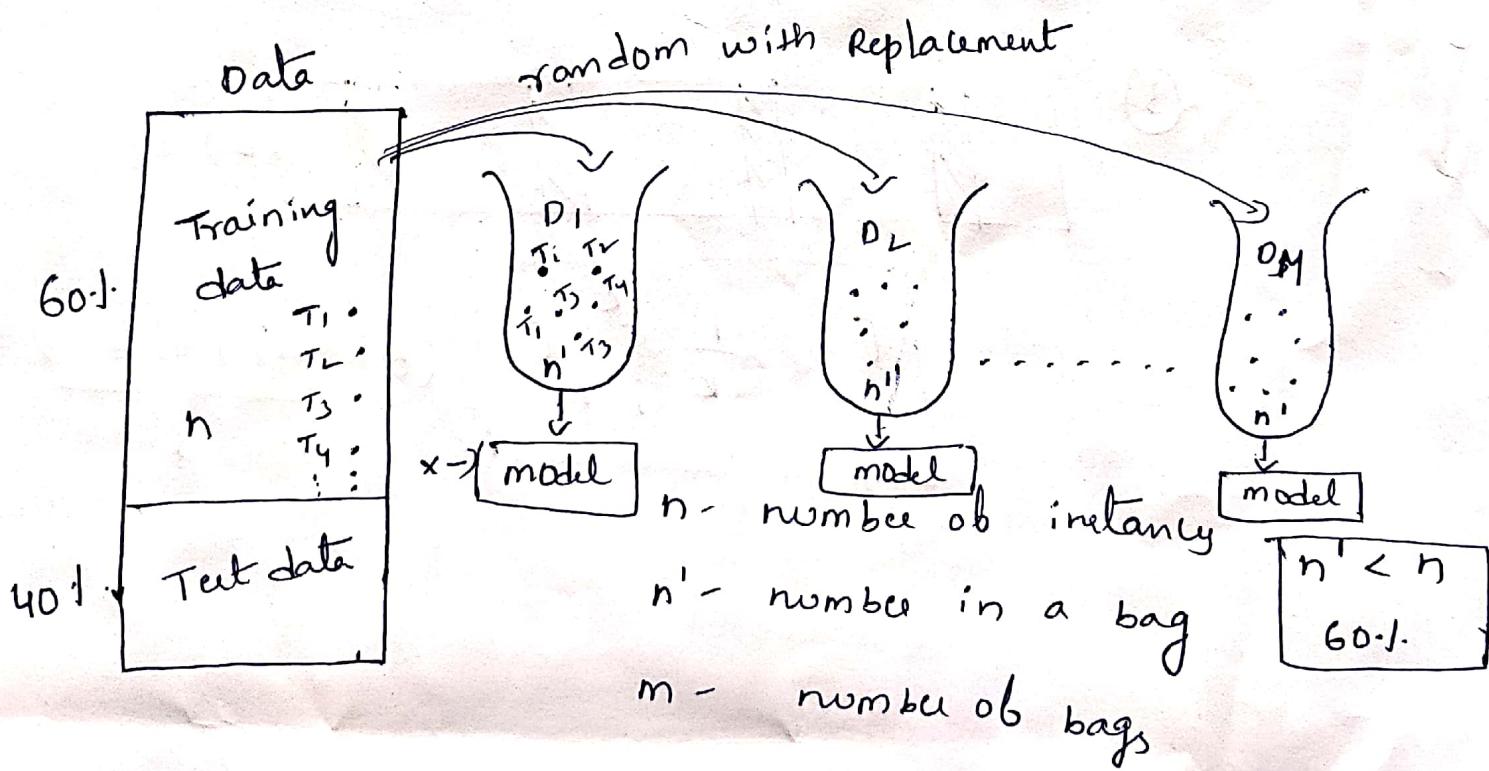
39

Bootstrap :- (8) Bootstrap aggregating - bagging

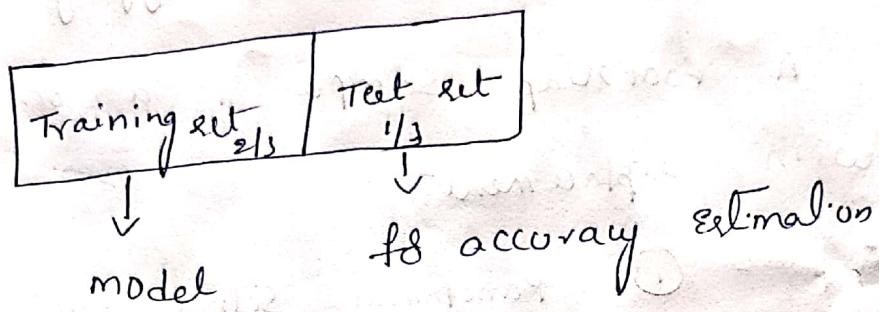
A Bootstrap Sample is a random sample conducted with replacement

- Steps:-
- ① Randomly select an observation from the original data
 - ② "write it down"
 - ③ "put it back" (i.e any observation can be selected more than once)

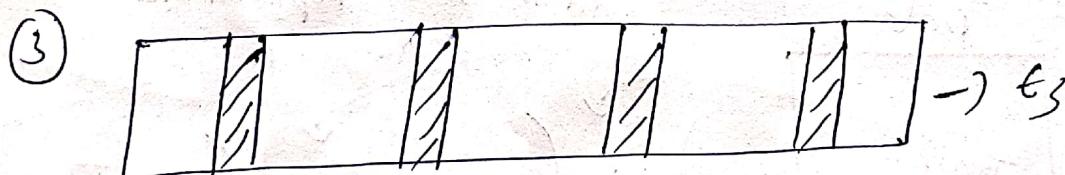
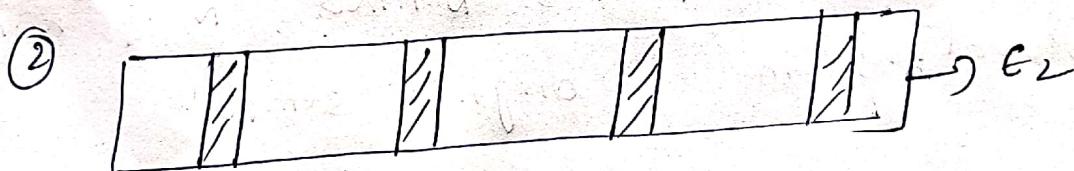
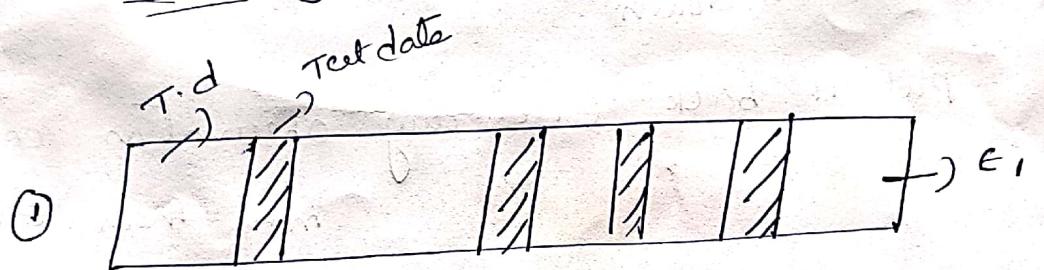
Repeat steps 1-3 n times, n is number of observations in the original sample.



Hold out method :-



Random subsampling (B) repeated Hold out method



$$F = \frac{1}{k} \sum_{i=1}^k e_i \quad (B) \quad \frac{e_1 + e_2 + e_3}{3}$$

Ada Boost.

- Given a set of d class-labeled tuples $(x_1, y_1) \dots (x_d, y_d)$.
- Initially all the weights of tuples are set the same ($1/d$).
- Generate k classifiers in K rounds. At round i
 - * Tuples from D are sampled to form a Training set D_i of the same size.
 - * Each tuple's chance of being selected is based on its weight.
 - * A classification model m_i is derived from D_i .
 - * Its error rate is calculated using D_i as a test set.
 - * If a tuple is misclassified, its weight is increased, otherwise it is decreased.
- $\text{Error}(m_i) = \sum_j^d w_j \times \text{err}(x_j)$
- weight of classifier m_i 's vote is

$$\log \frac{1 - \text{error}(m_i)}{\text{error}(m_i)}$$

Random forest

3.4

Another ensemble method called random forest. Each of the classifiers in the ensemble is a decision tree classifier. So that the collection of classifiers is a "forest".

Two methods to construct Random Forest

Forest - RI :- (Random input selection):

Randomly select, at each node, f attribute as candidates for the split at the node. The CART methodology is used to grow the trees to maximize size.

Forest - RC (random linear combinations) :-

Create new attribute (feature) that are a linear combination of existing attribute.