

### \* Classification :-

Classification is the task of learning a target function "g" that maps each attribute set "x" to one of pre-defined class labels "y".



The target function is known informally as a "classification model".

\* Descriptive Modelling :- A classification model can serve as an explanatory tool to distinguish the objects of diff classes.

→ for ex it would be useful for both biologist & officer to have a descriptive model that summarizes data shown.

Name	Body temp	Skin cover	Ciped BPrth	Aquatic creature	Aerial creature	Hair legs	Class label
Human	warm	hair	Yes	No	No	Yes	Mammal
python	cold	Scaler	No	No	No	No	Reptiles
frog	cold	none	No	Semi	No	Yes	Amphibian
whale	warm	hair	Yes	Yes	No	No	Mammal

\* Predictive Modelling :- A classification model can also be used to predict the class label of unknown records.

→ A classification model can be treated as a "black-box" that automatically assign a class label when presented with attribute set of an unknown record.

Name	Body temp	Skin cover	Ciped BPrth	Aquatic creature	Aerial creature	Hair legs	Class label
gila Monitor	cold	Scaler	No	No	No	Yes	Reptiles

We use a classification model built from data set to determine

\* general approach to solving a classification problem

A classification technique (or classifier) is a systematic approach to building classification model from an input set data.

Ex: decision tree classifier, rule-based classifier etc

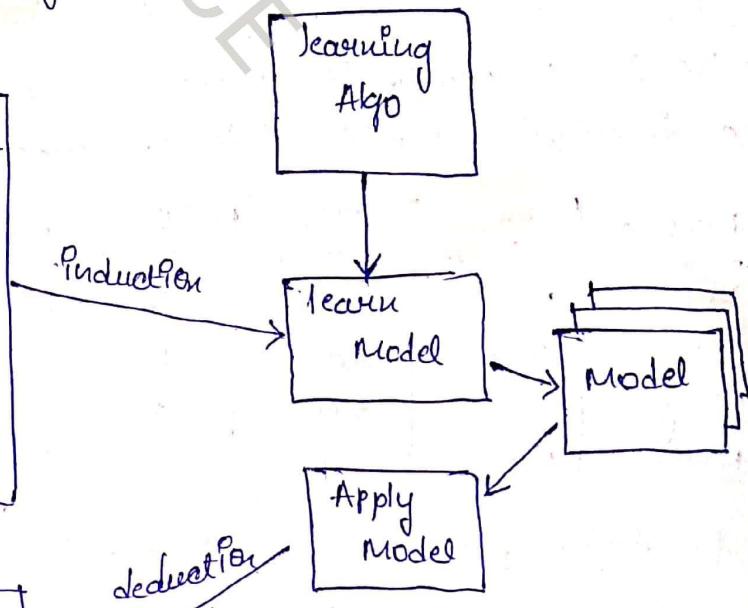
- Each technique employs a "learning Algo" to identify a model that best fits the relationships b/w the attribute set & class label of input data.
- the model generated by learning Algo should both fit the input data well & accurately predict the class label of records it has never before.

Tid	Attr1	Attr2	Class
1	Yes	Large	No
2	No	Small	No
3	Yes	Medium	No
4	No	Large	Yes
5	Yes	Large	No
6	Yes	Small	No

Training Set

Tid	Attr1	Attr2	Class
7	Yes	Large	?
8	No	Small	?
9	Yes	Small	?

Test Set



The evaluation of performance of a classification model is based on the counts of test records correctly & incorrectly predicted by the model.

These counts are tabulated in a table known as "confusion matrix"

		Predicted class	
		Class = 1	Class = 0
Actual class	Class = 1	$F_{11}$	$F_{10}$
	Class = 0	$F_{01}$	$F_{00}$

Table depicts the confusion Matrix for a binary classification problem, Entry  $F_{ij}$  in the table denotes the no of records from class  $i$  predicted to be of class  $j$ .

→ The performance Metric such as accuracy which is defined as follows

$$\text{Accuracy} = \frac{\text{No of correct predictions}}{\text{total no of predictions}} = \frac{F_{11} + F_{00}}{F_{10} + F_{11} + F_{01} + F_{00}}$$

The performance of a Model can be Expressed in terms of the Error-Rate which is given by

$$\text{Error-rate} = \frac{\text{No of wrong predictions}}{\text{total no of predictions}} = \frac{F_{10} + F_{01}}{F_{01} + F_{10} + F_{11} + F_{00}}$$

Most classification Algo seeks Model that "attain the highest accuracy or equivalently lowest Error Rate" when applied to test set.

\* decision tree Induction - Which is a simple yet widely used classification technique.

→ The tree has three types of nodes

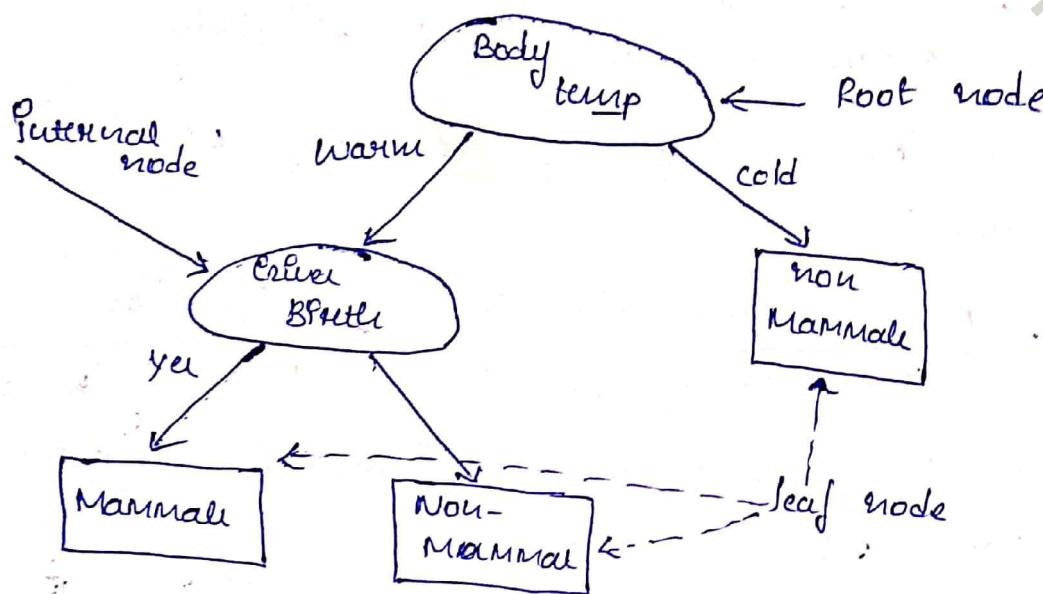
\* root node → that has no incoming Edge & zero or more outgoing edges

\* Internal node → Exactly one incoming Edge & two or more outgoing edges.

\* leaf node → nodes each of which has exactly one incoming edge

In a decision tree each leaf node is assigned to class label, the non-terminal nodes which include the root & other internal nodes contain attribute test conditions to separate records that have diff characteristics.

→ Consider the Ex as shown below



\* How to Build a decision tree e- One Such Algo Be "Hunter Algo" which Be the basis of Many Existing decision tree Production Algo Including ID3, C4.5, & CART

\* Hunter Algo e- In Hunter Algo, a decision tree is grown in a recursive fashion by partitioning the training records into successively purer subsets

→ Let  $D_t$  be the set of training records that are associated with node  $t$  &  $y = \{y_1, y_2, \dots, y_n\}$  be class labels the foll Be Recursive defn of Hunter Algo

\* Step 1 e- If all records in data belong to some class  $y_t$  then  $t$  is the leaf node labeled as  $y_t$

\* Step 2 e- If data contains records that belong to more than one class then choose a best attribute  $A$  to split the data

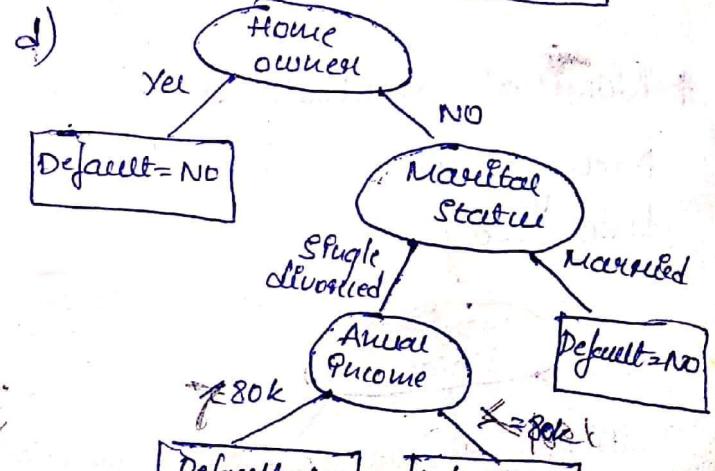
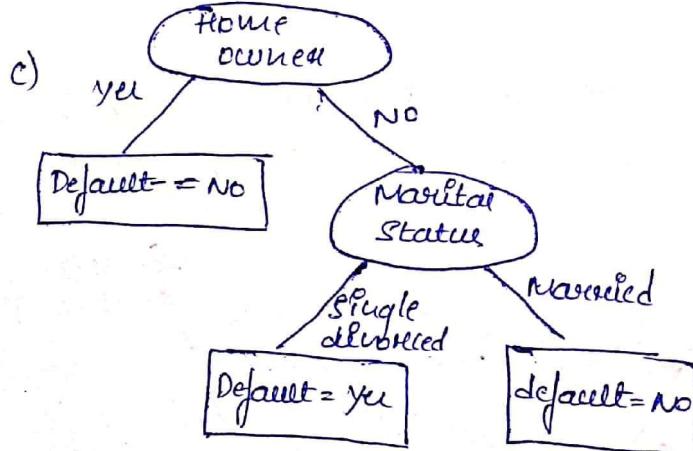
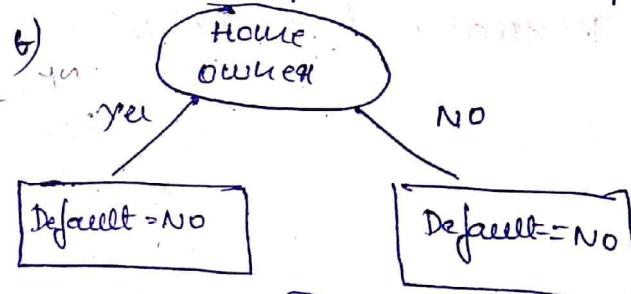
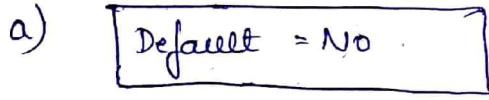
23

To partition the records into smaller subsets.

A child node is created for each outcome of test condition & records in Dt are distributed to children based on outcomes.

→ consider the foll Example problem of predicting whether loan applicant will repay the loan

TID	Home owner	Marital status	Annual Income	Defaulted Borrower
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	70k	Yes



Additional conditions are needed to handle the fall cases.

\* It will be possible for some of child nodes created in Step 2 to be empty.

→ i.e. there are no records associated with these nodes.

\* In Step 2 If all the records associated with Dt have identical attribute values then it is not possible to split them further.

\* design rules of decision tree Production-

A learning algo for producing decision tree must address the fall two rules.

\* How should the training records be split?

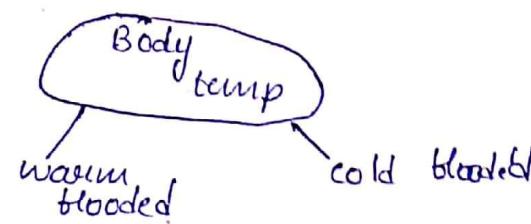
→ Attribute test condition

\* How should the splitting procedure stop?

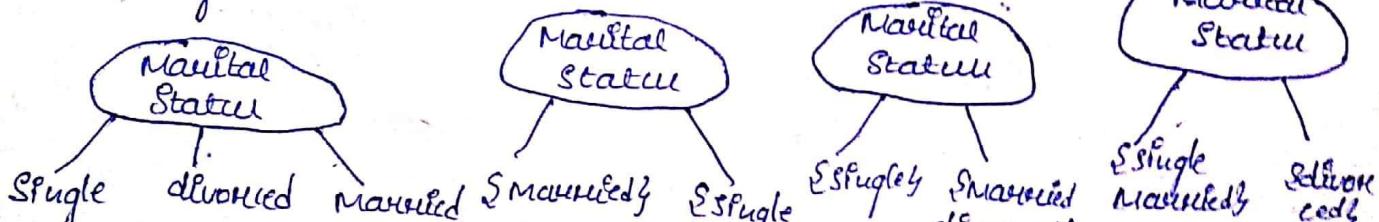
→ Stopping condition.

→ Method for Expressing Attr test condition

\* Binary attribute → the test condition for binary attr generates two potential outcomes

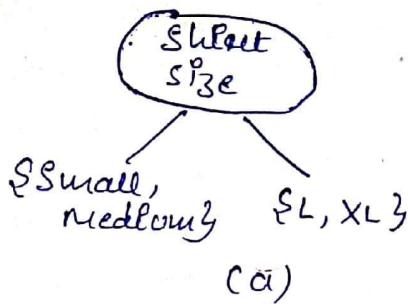


\* Nominal attribute → Since a Nominal attribute can have many values the test condition can be expressed in two ways.

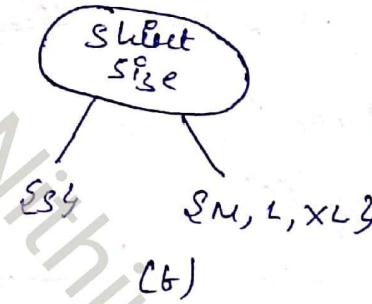


\* Ordinal attribute  $\rightarrow$  can also produce binary or Multiple Splits. (4)

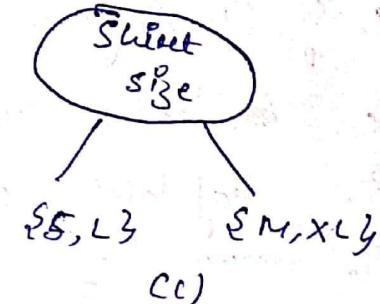
Ordinal attribute values can be grouped as long as the grouping does not violate the order property of attribute values.



(a)

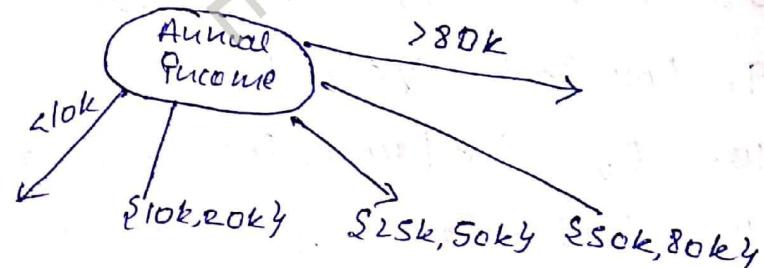


(b)



(c)

\* Continuous attribute  $\rightarrow$  For continuous attribute, the best condition can be expressed as comparison test ( $A < V$ ) OR ( $A \geq V$ ) with binary outcomes



Measures for Selecting the Best Split - Measures developed for selecting the best split are often based on "degree of purity" of child nodes

$\rightarrow$  The smaller the degree of purity, the more skewed the class distribution

$$* \text{Entropy } C(t) = - \sum_{i=0}^{C-1} p(i/t) \cdot \log_2 p(i/t)$$

$$* \text{Gini } C(t) = 1 - \sum_{i=0}^{C-1} [p(i/t)]^2$$

$$* \text{Classification Error} = 1 - \max_t [p_e(i/t)]$$

## \* Algorithm for decision tree Production

the Input to this Algo consists of training records "E" & attribute set "F"

### Tree Growth ( $E, F$ )

1. If Stopping-cond ( $E, F$ ) = true then
2. Leaf = Create Node ()
3. Leaf.Label = classify ( $E$ )
4. return Leaf
5. Else
6. Root = Create Node ()
7. Root.test-cond = Find\_Best\_Split ( $E, F$ )
8. Set  $V = \{v | v \text{ is a possible outcome of } Root.\text{test\_cond}\}$
9. For Each  $v \in V$  do
10.  $E_v = \{e | Root.\text{test\_cond}(e) = v \wedge e \in E\}$
11. child = Tree Growth ( $E_v, F$ )
12. add child as descendant of root & label the edge as  $v$
13. End For
14. End IF
15. Return Root.

After building the decision tree "Tree-pruning" Step can be performed to reduce the size of decision tree.

→ decision tree that are too large are susceptible to phenomenon known as "overfitting"

## \* Characteristics of decision tree Production

- \* It is non-parametric Approach for building classification Model.

\* decision tree are computationally inexpensive Making 2 (5)  
It possible to quickly construct even when training set  
size be very large.

\* decision trees. Especially smaller-sized tree are relatively  
easy to interpret.

\* decision tree Algo are quite robust to presence of  
noise.

\* The presence of redundant Atts. does not adversely affect  
the accuracy of decision tree.

\* Since decision tree use "top-down" recursive approach  
the no of records become smaller.

\* "Sub tree" can be replicated many times.

\* "test condition" involve using only one attribute at a  
time.

\* Model overfitting:- The Error committed by a classifier  
-on Model are generally divided into two types.

\* Training Error.

\* generalization Error.

\* "Training Error" also known as "Resubstitution Error" or  
"Apparent Error" is the no of misclassification error  
committed on training records.

\* "generalization Error" is the expected error of the model  
on previously unseen records.

A good Model Must have low training Error as well as  
low generalization Error.

\* Boot Strap- Boot Strap approach, the training records are sampled with replacement, i.e. a record already chosen from training is put back into the original pool of records so that it is equally likely to be drawn again.

→ A boot strap sample size  $N$  contains about 63.2% of the original data.

This approximation follows from the fact that probability a record chosen by a boot strap sample is

$$1 - (1 - 1/N)^N$$

$$1 - e^{-1} = 0.632$$

The sampling procedure repeated  $b$  times to generate  $b$  boot strap samples.

\* Method for comparing classifiers

\* Estimating a confidence interval for accuracy-

To determine the confidence interval, we need to establish the probability distribution that governs the accuracy measure.

→ Exp be a list of characteristic of binomial Exp

+ Exp consists of  $N$  independent trials, where trials have two possible outcomes → Success  
→ failure

+ the probability of Success  $p$ , in each trial is constant.

→ An Ex of a binomial Exp for counting no of heads that turn up when coin flipped  $N$  times

$$P(X=v) = \binom{N}{v} p^v (1-p)^{N-v}$$

→ Based on Normal distribution, the fold confidence Intervals for acc can be derived

$$P \left( -z_{\alpha/2} \leq \frac{\text{acc} - p}{\sqrt{p(1-p)/N}} \leq z_{1-\alpha/2} \right) = 1 - \alpha$$

\* Comparing the performance of two Model :- Consider a pair of Models  $M_1$  &  $M_2$  that are evaluated on two independent test Sets  $D_1$  &  $D_2$

→ let  $n_1$  denote no of records in  $D_1$ , &  $n_2$  denote no of records in  $D_2$

Suppose the Error rate for  $M_1$  on  $D_1$  is  $e_1$ ,  
the Error rate for  $M_2$  on  $D_2$  is  $e_2$

→ therefore the Variance of  $d$  can be computed as follows

$$\hat{\sigma}_d^2 \approx \hat{\sigma}^2 = \frac{e_1(1-e_1)}{n_1} + \frac{e_2(1-e_2)}{n_2}$$

where  $e_1(1-e_1)/n_1$ , &  $e_2(1-e_2)/n_2$  are the Variance of Error rates finally at  $(1-\alpha)$  % confidence level.

\* Comparing performance of two classifiers :- Suppose we want to compare the performance of two classifiers using k-fold cross-validation approach.

→ Initially data set  $D$  be divided into  $k$  Equal sized partitions, we then apply each classifier to construct a model from  $k-1$  of partition & test it on remaining partition. The step repeated  $k$  times.

$$\hat{\sigma}_{CV}^2 = \frac{\sum_{j=1}^k (d_j - \bar{d})^2}{k(k-1)}$$

where  $\bar{d}$  is average difference, for this approach we need to use a t-distribution to compute the confidence interval for  $d_t^{cv}$

$$d_t^{cv} = \bar{d} \pm t_{(1-\alpha), k-1} \sigma_d^{cv} . f$$

\* Rule Based classifier :-

A rule-based classifier is a technique for classifying records using a collection of "if... then..." rules.

- The rules for the Model are represented in a disjunctive normal form  $R = (H_1, V_1) \vee (H_2, V_2) \vee \dots \vee (H_k, V_k)$  where  $R$  is known as "rule Set" & " $H_i$ " are classification rules or disjunctions.
- Example of rule Set for vertebrate classification problem

$H_1$  : (Cervical\_Birth = no)  $\wedge$  (Aerial\_Creature = yes)  $\rightarrow$  Bird

$H_2$  : (Cervical\_Birth = no)  $\wedge$  (Aquatic creature = yes)  $\rightarrow$  Fish

$H_3$  : (Cervical\_Birth = yes)  $\wedge$  (Body temp = warm)  $\rightarrow$  Mammal

$H_4$  : (Cervical\_Birth = no)  $\wedge$  (Aerial\_Creature = no)  $\rightarrow$  Reptile

- Each classification rule can be Expressed in foll way

$H_i$  : (Condition<sub>p</sub>)  $\rightarrow$   $X_i$

The left-hand side of the rule is called "rule antecedent" or "precondition" It contains a conjunction of attribute tests

Condition<sub>p</sub> =  $(A_1, op_1 V_1) \wedge (A_2, op_2 V_2) \wedge \dots \wedge (A_k, op_k V_k)$

where  $(A_j; V_j)$  is an attribute-value pair & "op" is a logical operator chosen from set  $\{=, \neq, <, >, \leq, \geq\}$

Each attribute test  $(A_j, op_j V_j)$  is known as "conjunction".

- The right-hand side of the rule is called the "rule consequence" which contains one test of the form

→ A Rule R covers a record  $x$  if the precondition matches the attribute of  $x$ .

R is also said to be fired or triggered whenever it covers a given record.

→ The quality of a classification rule can be evaluated using measures such as

+ coverage

+ accuracy

\* The coverage of rule R is defined as fraction of records in D that trigger rule R.

$$\text{Coverage}(R) = \frac{|A|}{|D|} \rightarrow \text{the fraction of records that satisfy the antecedent}$$

\* the accuracy or confidence factor is defined as the fraction of records triggered by R whose class labels are equal to y.

$$\text{Accuracy}(R) = \frac{|Anyl|}{|A|} \text{ fraction of records that satisfy both antecedent \& consequent}$$

where  $|A| \rightarrow$  no of records that satisfy rule antecedent

$|Anyl| \rightarrow$  no of records that satisfy both antecedent & consequent

$|D| \rightarrow$  total no of records.

\* How a Rule-based classifier works-

A Rule-based classifier classifies a test record based on rules triggered by the record.

→ Consider the rule set for all vertebrates.

Name	Body temp	Skin cover	Cloacal birth	Aquatic creature	Aerial creature	Fan stage	Hibernation
lemon	warm	Fur	yes	No	No	yes	yes
turtle	cold	Scaler	no	Semi	no	yes	no
dogfish	cold	Scaler	yes	yes	No	No	No
shark	cold	Scaler	yes	yes	No	No	No

\* the first vertebrate which is not warm is lemon.

to be classified as a Mammal.

2(1)

\* the second vertebrate, which is turtle, triggers rules R<sub>1</sub> & R<sub>2</sub>, these conflicting classes must be resolved.

\* None of the rules are applicable to dogfish shark.

→ Mutually Exclusive rules :- the rules in rule set R are mutually exclusive if no two rules in R are triggered by same record.

This property ensures that every record is covered by at most one rule in R.

→ Exhaustive rules :- A rule set R has exhaustive coverage if there is a rule for each combination of attribute values.

This property ensures that every record is covered by at least one rule in R.

Assuming that "Body temp" & "Cervical Birth" are binary variables

→ Example of Mutually Exclusive & Exhaustive rule set

R<sub>1</sub> : (Body temp = Cold) → Non-Mammal

R<sub>2</sub> : (Body temp = Warm) ∧ (Cervical Birth = Yes) → Mammal

R<sub>3</sub> : (Body temp = Warm) ∧ (Cervical Birth = No) → Non-Mammal

→ Ordered rules :- By this approach, the rules in rule set are ordered in decreasing order of their priority, which can be defined in many ways (coverage, accuracy etc)

\* This avoids the problem of having conflicting classes predicted by multiple association rules.

→ Unordered rules :- this approach allows a test record to trigger multiple classification rules & consider the outcome of each rule as a vote for particular class.

The record is usually assigned to class that receives the

## \* Direct Method for Rule Extraction

"Sequential covering"

Algorithm is often used to Extract rules directly from data.

Rules are grown in a greedy fashion based on certain Evaluation Measure.

→ Sequential covering Algo is as shown

1. Let  $E$  be the training records &  $A$  be set of attribute value pairs  $\{A_j, V_j\}$
2. Let  $y_0$  be an ordered set of classes  $\{y_1, y_2, \dots, y_k\}$
3. Let  $R = \emptyset$  be initial rule. Let  $s$
4. For Each class  $y \in y_0 - \{y_k\}$  do
5. While Stopping condition is not met do
6.  $\leftarrow$  Learn-one-rule  $(E, A, y)$ .
7. Remove training records from  $E$  that covered by  $r$
8. Add  $r$  to the bottom of rule. Let  $s \rightarrow R \cup r$
9. End while
10. End for
11. Insert the default rule  $\emptyset \rightarrow y_k$  to the bottom of rule. Let  $R$ .

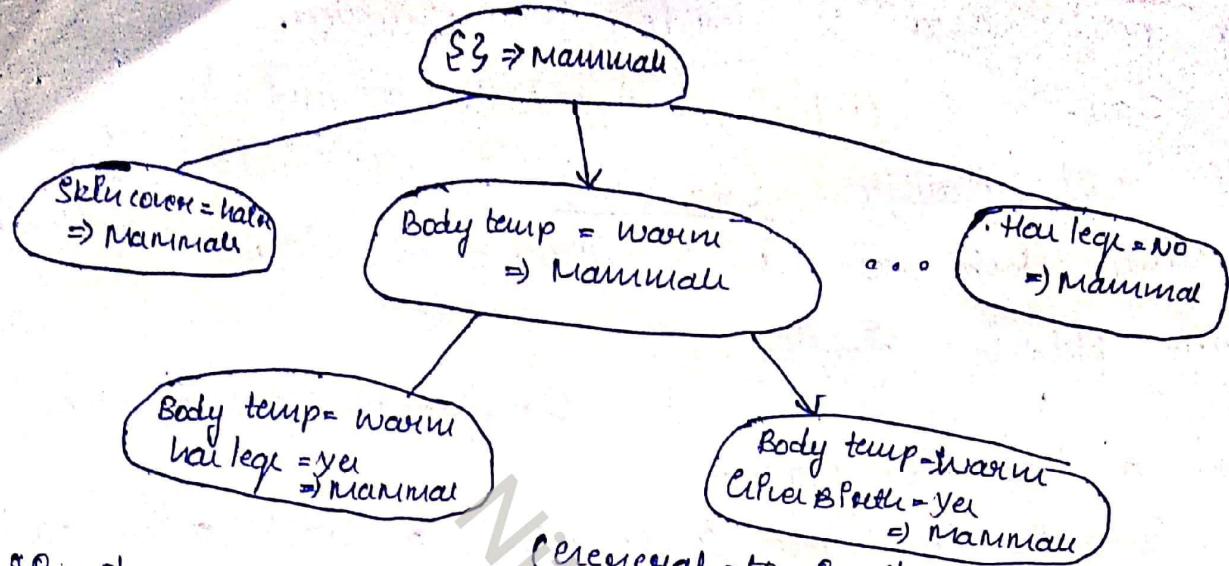
→ Rule-Growing Strategy: there are two common strategy

- see for growing a classification

\* general-to-specific

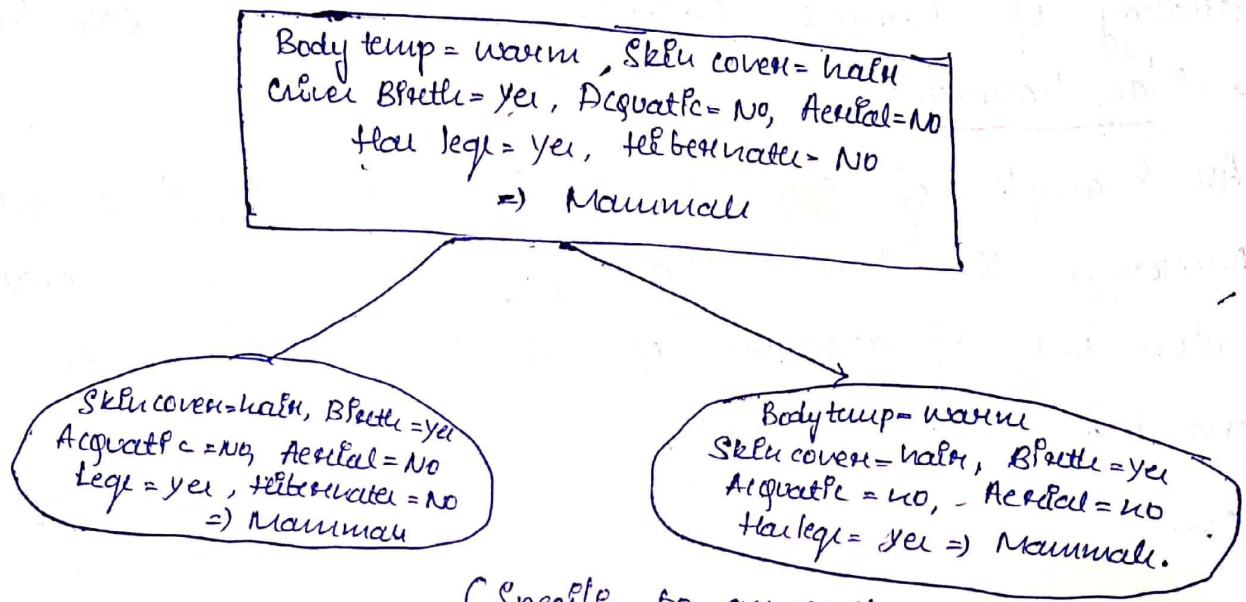
\* specific-to-general

\* "general-to-specific" Strategy in which an initial rule  $\emptyset \rightarrow y$  is created, where the left-hand side is an empty set & right-hand side contains the target class. The rule has poor quality because it covers all  $\text{ex}$  in



- \* "SpecifPc-to-general" Strategy, One of the positive Examples, is randomly chosen as RulePc Seed for rule-growing process.

during the refinement step, the rule is generalized by removing one of the conjuncts so that it can cover more positive Examples.



(Specific-to-general)

- \* characteristic of Rule-based classifier

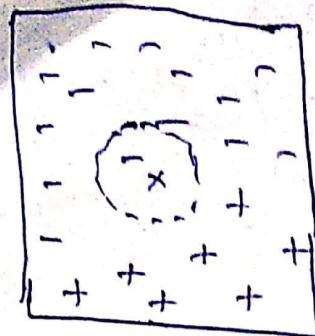
\* the Expressiveness of classifier is almost Equivalent to that of a decision tree

Because a decision tree can be represented by a set of Mutually Exclusive & Exhaustive Rules.

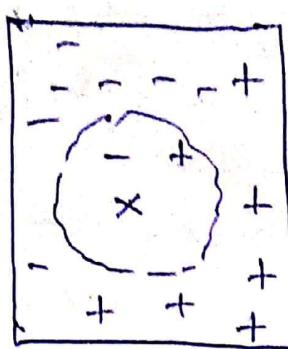
- \* Rule-based classifiers are generally used to produce descriptive models and explanatory models.

- able performance to deal with tree classifiers.
- \* the class-based ordering Approach adopted by Many Rule-based classifiers (CRIPPER) is well suited for handling data sets with imbalanced class distribution.
- \* Nearest-Neighbour classifier is part of the classification framework involving a two-step process
  - \* an Inductive Step for constructing a classification Model from data.
  - \* a deductive Step for applying Model to test Examples
- "decision tree classifiers" & "rule-based classifiers" are good Example of "Eager learners"  
And meanwhile the techniques that employ the strategy of Nearest-Neighbour classification are known as "Lazy learners"
- An Example of Lazy learners is "Rule classifier" which memorizes the entire training data & performs classification only if attribute of test instance matches one of training example exactly.

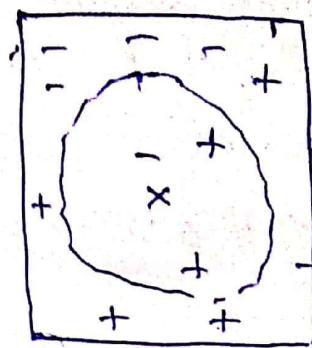
- ~~- If we int k is specified along a new sample
  - Compute the distance b/w ilp sample & sample training data by using  $(P_i - Q_j)^2$
  - Euclidian distance
  - Apply decision boundary
  - Select k entries in database which are closest to sample
  - Then this be classified to new sample~~



1-nearest neighbor



2-nearest neighbor



3-nearest neighbor

→ One way to make this approach more flexible & to find all training ex that are relatively similar to the attribute of test example, which are known as "nearest neighbour". can be used to determine the class label of test example.

→ the 1-, 2-, 3-nearest neighbour of data point located at the center of each circle, the data point be classified based on the class label of its neighbours.

In the case where the neighbours have more than one label, the data point be assigned to majority class of its nearest neighbour.

\* Algorithm 6- the Algo compute the distance bw similarity bw Each test Example  $\mathbf{z} = (x^1, y^1)$  & all the training Examples  $(x, y) \in D$  to determine the nearest-neighbor  $\mathbf{z}$ .

1. Set  $k$  be no of nearest neighbour &  $D$  be set of training ex
2. For Each test Example  $\mathbf{z} = (x^1, y^1)$  do
3. compute  $d(x^1, x)$ , the distance bw  $\mathbf{z}$  & Every ex  $(x, y) \in D$
4. Select  $D_2 \subseteq D$ , the set of  $k$  closest training ex of  $\mathbf{z}$

$$5. y' = \arg \max_{y'} \sum_{(x_i, y_i) \in D_2} I(V = x_i)$$

End for

→ Once the Nearest-Neighbour set is obtained, the test  $\underline{x}_k$  is classified based on the Majority class of the nearest neighbour.

\* Majority voting  $y' = \operatorname{argmax} \sum_{(x_i, y_i) \in D_2} I(V = y_i)$

→ Using the distance-weighted voting scheme the class label can be determined as follows:

\* Distance-weighted voting  $y' = \operatorname{argmax} \sum_{(x_i, y_i) \in D_2} w_i \cdot I(V = y_i)$

#### \* Characteristics of Nearest Neighbour classifier

\* Nearest-Neighbour classification is part of a more general technique known as "Pertaining-based learning" which uses specific training perturbations to make prediction without having to maintain an abstraction derived from data.

\* Lazy learners such as nearest-neighbour classifiers do not require model building.

\* Nearest-Neighbour classifiers make their prediction based on local information.

\* Nearest-Neighbour classifiers can produce anisotropy shaped decision boundaries.

The decision boundaries of nearest-neighbour classifiers have "high variability" because they depend on the composition of training examples.

\* Nearest-Neighbour classifiers can produce wrong prediction unless the appropriate proximity measure & data preprocessing steps are taken.

\* Bayesian Classification

\* Bayes theorem is a Statistical principle for combining prior knowledge of class with new evidence gathered from data.

→ let  $x \in y$  be a pair of Random Variables, Then the joint probability  $p(x=x, y=y)$  refers to the probability that variable  $x$  will take on the value  $x$  & Variable  $y$  take on value  $y$ .

The joint & conditional probabilities for  $x \in y$  are related in following way

$$p(x,y) = p(y|x) \times p(x) = p(x|y) \times p(y)$$

Rearranging the expression leads to formula known as "Bayes theorem"

$$p(y|x) = \frac{p(x|y) p(y)}{p(x)}$$

\* Using Bayes theorem for classification :- Let  $x$  denote the attribute set &  $y$  denote the class variable, If the class variable has non-deterministic relationship with attribute then we can treat  $x \in y$  as Random Variable & capture their relationship using  $p(y|x)$ .

The conditional probability is also known as "posterior probability for y", as opposed to the "prior probability p(y)"

→ during the training phase, we need to learn the posterior probability  $p(y|x)$  for every combination of  $x$  &  $y$  based on info gathered in training data.

By knowing these probabilities, a test record  $x'$  can be

$$-p(y'|x')$$

→ Let's consider the full training set.

Tid	Home owner	Marital status	Annual Income	Borrower
1	yes	Single	125k	No
2	No	Married	100k	No
3	Yes	Single	85k	No
4	No	Divorced	60k	No
5	No	Single	120k	Yes
6	No	Single	100k	No
7			:	:
8	No	Married	80k	Yes

Suppose we are given a test record with full attribute set  $x = (\text{Home owner} = \text{No}, \text{Marital status} = \text{Married}, \text{A. Income} = 120k)$ . To classify the record we need to compute posterior probabilities  $p(\text{Yes}|x)$  &  $p(\text{No}|x)$  based on the info available in training data.

If  $p(\text{Yes}|x) > p(\text{No}|x)$  then the record is classified as "Yes" otherwise "No".

→ Bayes theorem is useful because it allows us to express posterior probability in terms of prior probability  $p(x)$ .

The "class-conditional" probability  $p(x|y)$  is evidence  $p(x)$

$$p(y|x) = \frac{p(x|y)p(x)}{p(x)}$$

When comparing the posterior probability for different values of  $x$ , the denominator term  $p(x)$  is always constant & thus can be ignored.

- \* Naive Bayes' classifier :- Estimate the class-conditional probability by assuming that the attributes are conditionally independent, given the class label  $X$ .
- The conditional independence assumption can be formally stated as follows

$$P(X|Y=y) = \prod_{i=1}^d P(X_i|Y=y)$$

where each attribute set  $X = \{X_1, X_2, \dots, X_d\}$  consists of  $d$  attributes.

- \* Conditional Independence :- Let  $X, Y, Z$  denote three sets of random variables, the variables in  $X$  are said to be conditionally independent of  $Y$ , given  $Z$ , if the following condition holds

$$P(X|Y, Z) = P(X|Z).$$

#### → Characteristics of Naive Bayes Classifier

- \* they are robust to isolated noise points because such points are averaged out when estimating conditional probabilities from data.
- \* they are robust to irrelevant attributes.
- \* correlated attributes can degrade the performance of classifier because the conditional independence no longer holds for such attributes.

#### \* Bayesian Belief Networks (BBN) :-

provides a graphical representation

- representation of the probabilistic relationships among a set of random variables.

There are two key elements of Bayesian Net

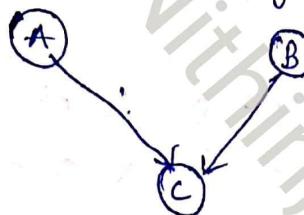
- \* A directed Acyclic graph (DAG) Encoding the dependencies relationships among a set of variables.

\* A probability table allocating each node to its parent.

- Parent node.

→ Consider three random variables A, B, C in which A & B are independent. Variable C each has a direct influence on third variable.

The relationship among the variables can be summarized into the directed acyclic graph as shown.



e.g.: Represent probabilistic relationship using directed acyclic graph.

→ conditional independence: - A node in a BBN is conditionally independent of its non-descendants, if its parents are known.

\* If a node X does not have any parents, then the table contains only prior probability  $p(x)$ .

\* If a node X has only one parent Y, then the table contains conditional probability  $p(x|y)$ .

\* If a node X has multiple parents  $\{Y_1, Y_2, \dots, Y_k\}$  then table contains conditional probability  $p(x|Y_1, Y_2, \dots, Y_k)$ .

→ Algorithm

1. Let  $T = (x_1, x_2, \dots, x_d)$  denote a total order of variables

2. for  $j=1$  to  $d$  do

3. let  $X_{T(j)}$  denotes the  $j^{\text{th}}$  highest order variable in  $T$

4. let  $\pi(X_{T(j)}) = \{x_{T(1)}, x_{T(2)}, \dots, x_{T(j-1)}\}$  denote set of variables preceding  $X_{T(j)}$

5. Remove variables from  $\pi(X_{T(j)})$  that don't affect  $X_j$

6. Create all the factors  $x_{T(j)}$  for  $j = 1, 2, \dots, d$

## → Characteristics of BBN

- \* BBN provide an approach for capturing prior knowledge of particular domain using graphical Model.
- \* Construction of N/w can be time consuming & require a large amount of Effort.
- \* BBN are well suited to dealing with incomplete data.
- \* Because the data be combined probabilistically with prior knowledge, the Method be quite robust to Model overfitting.