

Unit - 3

Classification

Classification is a data mining function that assign items in a collection to target categories or classes.

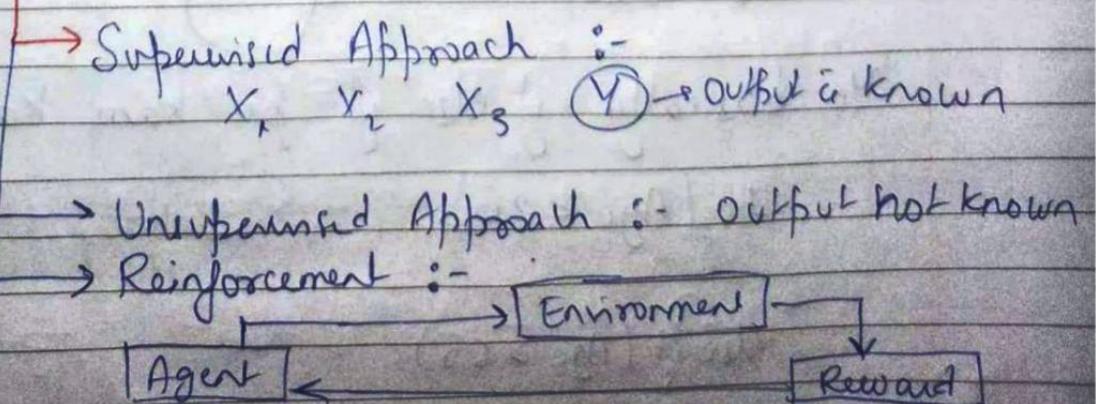
The goal of classification is to accurately predict the target class for each class in the data.

Example of cases where the data analysis task is classification :-

(a) A bank officer want to analyze the data in order to know which customer are risky or which are safe.

(b) A marketing manager at a company needs to analyze a customer with a given profile, who will buy a new computer.

ML (Machine Learning)



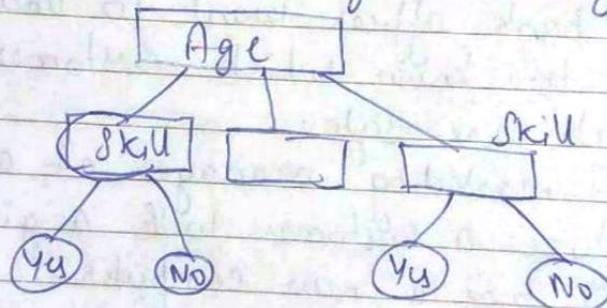
* Classification comes under Supervised approach

In classification, basically what we do is to learn understand the data and then by using known output we can predict the output of unknown data.

Decision Tree Induction

Decision Tree is a technique used in classification.

It will not be always a binary tree.



for attribute →



for result →



Ye basically if else if jpe kam kar shah.

If $\{Age < 20\}$

else if $\{Age > 60\}$

else ?

Now, whenever we choose different attributes or top decision tree will be different.

So, we need to find which attribute will be well suited for top. This process is known as Attribute Selection.

For this we measure Entropy.

$$\text{Entropy} = -\sum P_i \log_2(P_i) \text{ bit.}$$

<u>Q</u>	S.No	Age	Income	Class
1.	Youth	H		No
2.	Youth	H		No
3.	Middle age	H		Yes
4.	Senior	M		Yes
5.	Senior	L		Yes
6.	Senior	L		No
7.	Middle age	L		Yes
8.	Youth	M		No
9.	Youth	L		Yes
10.	Senior	M		Yes
11.	Youth	M		Yes
12.	Middle age	M		Yes
13.	Middle age	H		Yes
14	Senior	M		No

No. of clause = 2.
 Step 1 → Count No. of Yes & No.

$$\text{Yes} = 9 \\ \text{No} = 5$$

$$P(\text{Yes}) = \frac{9}{14} \quad P(\text{No}) = \frac{5}{14}$$

Step 2 → find information in data (Entropy)

$$\text{info}(D) = - \sum P_i \log_2(P_i)$$

$$= - \left(\frac{9}{14} \log_2 \left(\frac{9}{14} \right) + \left(\frac{5}{14} \times \log_2 \left(\frac{5}{14} \right) \right) \right)$$

$$= -(-0.4098 + (-0.5305))$$

$$= 0.940 \text{ bit}$$

Step 3:- find info of age :-

	Yes	No	Sum
Youth	2	3	5
Middle Age	4	0	4
Senior	3	2	5

$$\begin{aligned} \text{info(age)} &= - \left[\frac{5}{14} \times \left(\frac{2}{5} \log_2 \left(\frac{2}{5} \right) + \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right) \right. \\ &\quad + \left. \frac{4}{14} \times \left(\frac{4}{4} \times \log_2 \left(\frac{4}{4} \right) + 0 \times \log_2 0 \right) \right] \\ &\quad + \frac{5}{14} \times \left(\frac{3}{5} \log_2 \left(\frac{3}{5} \right) + \frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right) \end{aligned}$$

$$= -[0.347 + 0 + 0.347]$$

$$= 0.694 \text{ bits}$$

$$\text{Gain} = \text{info}(D) - \text{info}(\text{Age})$$

$$= 0.940 - 0.694$$

$$= 0.246 \text{ bits}$$

Step 4:- find info of income :-

	Yes	No	Sum
H	2	2	4
M	3	1	4
L	4	2	6

$$\text{info}(\text{income}) = -\left[\frac{4}{14} \times \left(\frac{2}{4} \log_2 \left(\frac{2}{4}\right) + \frac{2}{4} \log_2 \left(\frac{2}{4}\right)\right) + \right.$$

$$\left. \frac{4}{14} \times \left(\frac{3}{4} \log_2 \left(\frac{3}{4}\right) + \frac{1}{4} \log_2 \left(\frac{1}{4}\right)\right) + \right.$$

$$\left. \frac{6}{14} \times \left(\frac{4}{6} \log_2 \left(\frac{4}{6}\right) + \frac{2}{6} \log_2 \left(\frac{2}{6}\right)\right) \right].$$

$$= -[-0.2857 + 0.2326 + 0.3936]$$

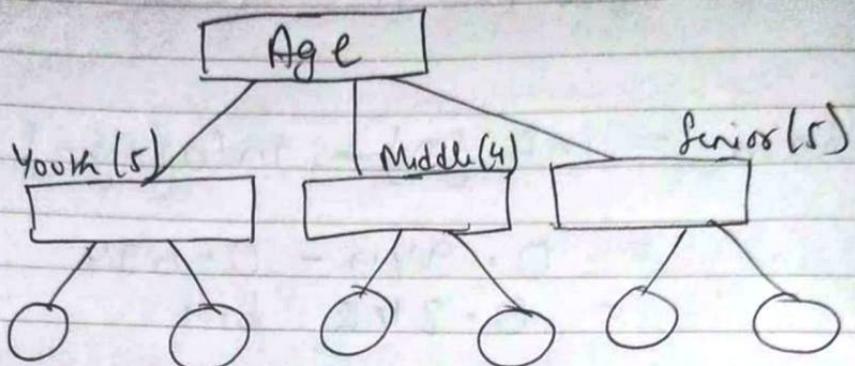
$$= 0.911 \text{ bits}$$

$$\text{Gain} = \text{info}(D) - \text{info}(\text{income})$$

$$= 0.940 - 0.911$$

$$= 0.029 \text{ bits}$$

Class of Age is higher than Class of income. So it will come on top.



Drawback :- 1) Bias :- it is favouring the attribute which has more no. of partitions.

Name of this Algo is ID3 :

Algorithm :-

Generate decision-tree (D, Attribute, Selection method)

- 1) Create a Node N
- 2) If D contains the tuples belonging to same class.
 - return N with Label of class
- 3) If attribute list is empty then
 - return N as leaf node with majority in D.
- 4) Apply attribute selection criteria to find the best split.
- 5) Attribute list = Attribute list - best selected attribute

6. For every j in the split

If D_j is empty.

Attack a leaf node with majority of class as labels.

Else

return Node_Genrate_decision-tree(D_j)

7. Return N.

$$* \text{Split_info}_A(D) = - \sum_{i=1}^v \left(\frac{d_i}{D} \right) \log_2 \left(\frac{d_i}{D} \right).$$

Split_info of Info(B)

$$= - \left(\frac{4}{14} \log_2 \frac{4}{14} + \frac{6}{14} \log_2 \frac{6}{14} + \frac{4}{14} \log_2 \frac{4}{14} \right).$$

$$= 1.55 \text{ bit}$$

$$\text{Gain Ratio} = \frac{\text{Info}(D)}{\text{Split_info}_A(D)}$$

$$= \frac{0.029}{1.55} = 0.018$$

* Select that whose gain ratio is high.

Gini Index :-

Gini Index measures that the impurity of data.

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

$$P(Y_1) = \frac{9}{14}$$

$$P(N_0) = \frac{5}{14}$$

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

$$= 0.46$$

= 46% impurity.

$$\text{Gini}_A(D) = \sum_{i=1}^v P_i / D \text{ Gini}(D)$$

* Gini is applied when split is Binary.

So, these 3 were the 3 techniques for designing Decision tree.

* Now the Question arises till what height we should do it.

We just need to draw till a certain level.
We will not draw after that level.

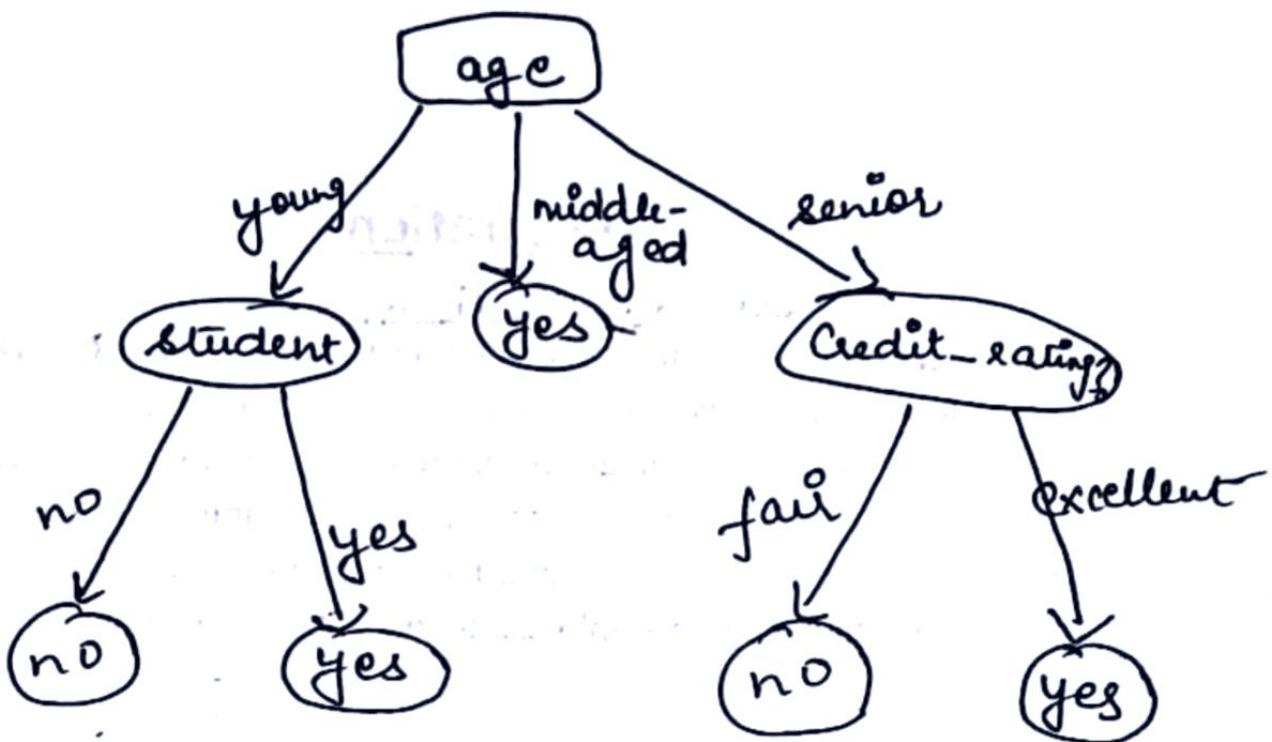
2 Ways to decide height

- (i) Pre planning
- (ii) Post planning.

Decision Tree induction

It is a method of learning the decision trees from the training sets. The dataset is broken down into smaller subsets and is present in the form of nodes of a tree. The tree structure has a root node, internal nodes or decision nodes, leaf nodes and branches. The root node is the topmost node.

- Each internal node denotes a test of an attribute, each branch denotes the outcome of a test, and each leaf node holds a class label.
The topmost node in the tree is the root node.



Decision tree induction algorithm

developed a decision tree algorithm known as ID3 (Iterative Dichotomiser) later, he presented C4.5, which was the successor of ID3. ID3 and C4.5 adopt a greedy approach. In this algo, there is no backtracking; the trees are constructed in a top-down recursive divide and conquer manner.

Tree Pruning

Tree Pruning is performed in order to remove anomalies in the training data due to noise or outliers. The pruned trees are smaller and less complex.

Tree - Pruning Approaches

There are two approaches to prune a tree -

Pre-pruning → The tree is pruned by halting its construction early.

Post-pruning → This approach removes a sub-tree from a fully grown tree.

Key factors

Entropy ↴

It refers to a common way to measure impurity. In the decision trees, it measures the randomness or impurity in data sets.

Information Gain ↴

Information gain refers to the decline in entropy after the dataset is split.

It is also called entropy reduction.

Building a decision tree is all about discovering attributes that return the highest data gain.

Note - Bayes Theorem

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Day	Outlook	Temp	Humidity	Wind	Play Cricket?
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Q → Find the probability to play cricket on 15th day when conditions are —

Temp = cool Humidity = High
 Wind = strong Outlook = Sunny

Step-1

Total(14)

Yes(9)

No(5)

$$P(\text{Yes}) = \frac{9}{14}$$

$$P(\text{No}) = \frac{5}{14}$$

expected %
Total no. of
outcome

Step-2

Wind(14)

weak(8)

strong(6)

Yes(6)

No(2)

Yes(3)

No(3)

$$P(\text{weak} | \text{Yes}) = \frac{6}{9}$$

$$P(\text{Strong} | \text{Yes}) = \frac{3}{9}$$

$$P(\text{weak} | \text{No}) = \frac{2}{5}$$

$$P(\text{Strong} | \text{No}) = \frac{3}{5}$$

Step-3

Humidity(14)

High(7)

Normal(7)

Yes(3) No(4)

Yes(6) No(1)

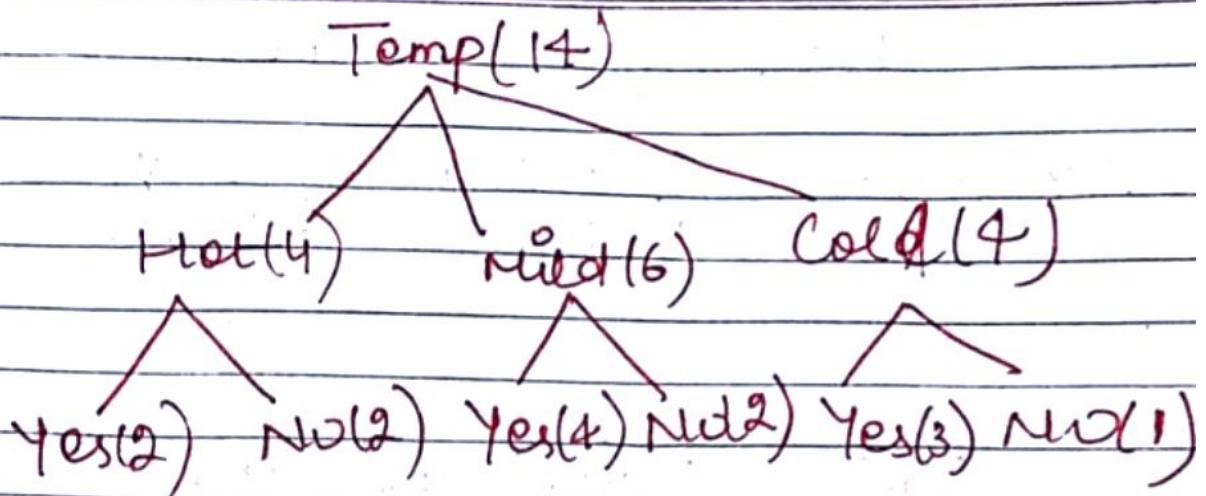
$$P(\text{High} | \text{Yes}) = \frac{3}{7}$$

$$P(\text{Normal} | \text{Yes}) = \frac{6}{7}$$

$$P(\text{High} | \text{No}) = \frac{4}{5}$$

$$P(\text{Normal} | \text{No}) = \frac{1}{5}$$

Step 4



$$P(\text{Hot} | \text{Yes}) = \frac{2}{9}$$

$$P(\text{Mild} | \text{Yes}) = \frac{4}{9}$$

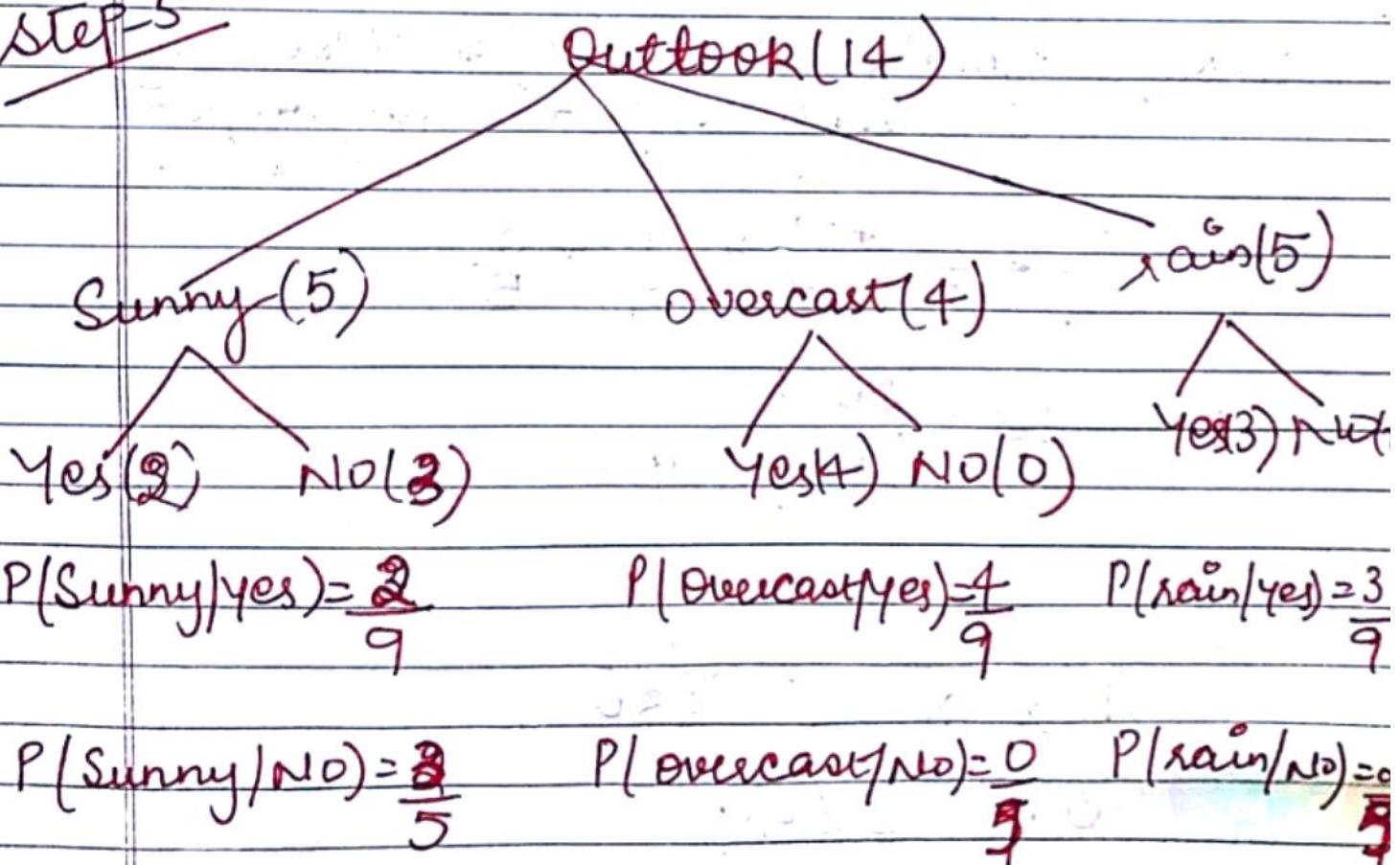
$$P(\text{Cool} | \text{Yes}) = \frac{3}{9}$$

$$P(\text{Hot} | \text{No}) = \frac{2}{5}$$

$$P(\text{Mild} | \text{No}) = \frac{2}{5}$$

$$P(\text{Cool} | \text{No}) = \frac{1}{5}$$

Step 5



$$P(\text{Sunny} | \text{Yes}) = \frac{2}{9}$$

$$P(\text{Overcast} | \text{Yes}) = \frac{4}{9}$$

$$P(\text{Rain} | \text{Yes}) = \frac{3}{9}$$

$$P(\text{Sunny} | \text{No}) = \frac{3}{5}$$

$$P(\text{Overcast} | \text{No}) = \frac{0}{5}$$

$$P(\text{Rain} | \text{No}) = \frac{2}{5}$$

Let $X = \{ \text{sunny}, \text{cool}, \text{high}, \text{strong} \}$

$$\begin{aligned} P(X/\text{Yes}) &= P(\text{Yes}) * P(\text{sunny}/\text{Yes}) * P(\text{cool}/\text{Yes}) \\ &\quad * P(\text{High}/\text{Yes}) * P(\text{Strong}/\text{Yes}) \\ &= \frac{9}{14} * \frac{2}{9} * \frac{3}{9} * \frac{3}{9} * \frac{3}{9} \\ &= 0.00529 \end{aligned}$$

$$\begin{aligned} P(X/\text{No}) &= P(\text{No}) * P(\text{sunny}/\text{No}) * P(\text{cool}/\text{No}) \\ &\quad * P(\text{High}/\text{No}) * P(\text{Strong}/\text{No}) \end{aligned}$$

$$= \frac{5}{14} * \frac{3}{5} * \frac{1}{5} * \frac{1}{5} * \frac{3}{5}$$

$$= \frac{18}{7 \times 125} = 0.02057$$

Probability of No is high

so answer should be "No".

→ RULE BASED CLASSIFICATION:

Rule based classifier uses a set of IF THEN rules for classification

If {condition} Then {condition}

↳ Rule antecedent ↳ Rule consequent

For eg - Rid Age Income Student Comp -

1	youth	high	NO	no
2	youth	high	NO	no
3	medium	high	NO	yes
4	old	medium	NO	yes
5	old	low	yes	yes
6	old	low	yes	no
7	medium	low	yes	yes
8	youth	medium	NO	no
9	youth	low	yes	yes

10	old	medium	Yes	Yes
11	youth	medium	Yes	Yes
12	medium	medium	No	Yes
13	medium	high	Yes	Yes
14	old	medium	No	No

R1: If age = 'youth' & student = 'yes'
then buy-computer = yes

R2: If age = 'old' & income > 25000 &
student = No then buy-computer = No

Rule has —

- i) Coverage \rightarrow fraction of records that satisfy antecedent of a rule.
- ii) Accuracy \rightarrow fraction of records that satisfy both antecedent & consequent of a rule.

for ex — R1: If age = youth & student = 'yes'
then buy-computer = yes.

$$\therefore \text{coverage} = 2/14 = 14.28\%$$

$$\text{accuracy} = 2/2 = 100\%.$$

Characteristics of a Rule based classifier

If we convert the result to classification rules, then rules would be mutually exclusive & exhaustive at same time.

a) ^{exclusive} Mutually exhaustive rules $R_1 \cap R_2 = \emptyset$

Classifier contains mutually exclusive rules. If then rules are independent of each other.

— Every record is covered by atmost one rule.

For eg - R_1 - Age & Income
 R_2 - Rid & student

b) exhaustive - Classifier has exhaustive coverage if it accounts for every possible combination of attribute :
— Each record is covered by atleast one rule.

eg (id, student) (id, income)
(age, income) (id, age, income)

But this is not possible since a record may trigger more than one rule or record may not trigger any rule.

- If more than one rule is triggered, need conflict resolution.

i) Size ordering → Assign the highest priority to the triggering rules that has the toughest requirement i.e. with most attribute test.

For eg -

R_1 - 3 Attribute

R_2 - 5 Attribute

$\boxed{R_3 - 11 \text{ Attribute}}$ → chosen.

(ii) Class-based ordering — Rules that belong to same class.

$$C_1 >= 25$$

$$\boxed{C_2 \leq 25 - 30}$$

$$C_3 > 50$$

iii) Rule-based ordering — rules are organised in one long priority list according to some measure of rule quality i.e. by coverage,

Eg $R_1 : (\text{give-birth} = \text{no}) \wedge (\text{can fly} = \text{yes}) \rightarrow \text{Bird}$

$R_2 : (\text{give birth} = \text{no}) \wedge (\text{live in water} = \text{yes}) \rightarrow \text{fish}$

$R_3 : (\text{give-birth} = \text{yes}) \wedge (\text{Blood type} = \text{warm}) \rightarrow \text{Mammals}$

$R_4 : (\text{give-birth} = \text{no}) \wedge (\text{can fly} = \text{no}) \rightarrow \text{Reptiles}$

$R_5 : (\text{live in water} = \text{sometimes}) \rightarrow \text{Amphibians}$

<u>Eg</u>	Name	BT	GB	can fly	live in water	cls
	Hawk	warm	no	yes	no	?
	Bear	warm	yes	no	no	?

R_1 covers hawk class bird

R_3 covers bear class mammals.

<u>Eg</u> →	Name	BT	GB	Can fly	Live in water	cls
	Turtle	need	no	no	sometimes	?
	Shark	cold	yes	no	sometimes	?

Turtle triggers both R_1 & R_5

↳ ordering rule or voting

default shark triggers none of rules

↳ default rules

Approaches for rule based classifier

- i) Direct method → extract rules direct from data.
eg— sequential covering algorithm.
- ii) Indirect method → extract rules from other classifier.
eg → decision tree.

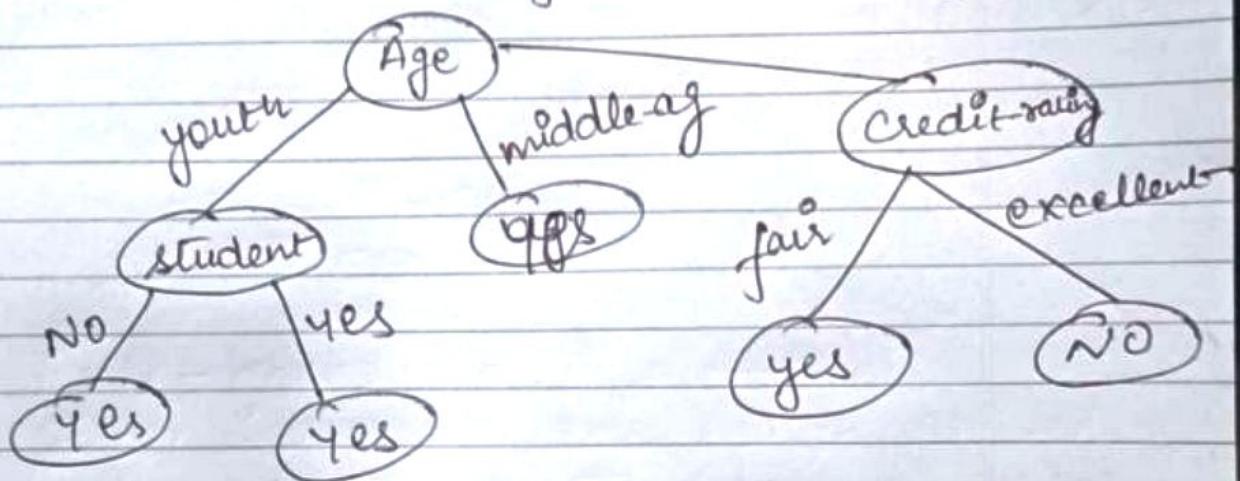
R1: If Age = youth AND student = no Then buy_computer = NO

R2: If Age = youth AND student = yes Then buy_computer = YES

R3: If Age = middle Then buy_computer = YES

R4: If Age = senior AND credit_rating = fair Then buy_computer = YES

R5: If Age = senior AND credit_rating = excellent Then buy_computer = NO



Methods for estimating a classifier's accuracy—

- i) Holdout method, random subsampling
- ii) Cross-validation
- iii) Bootstrap

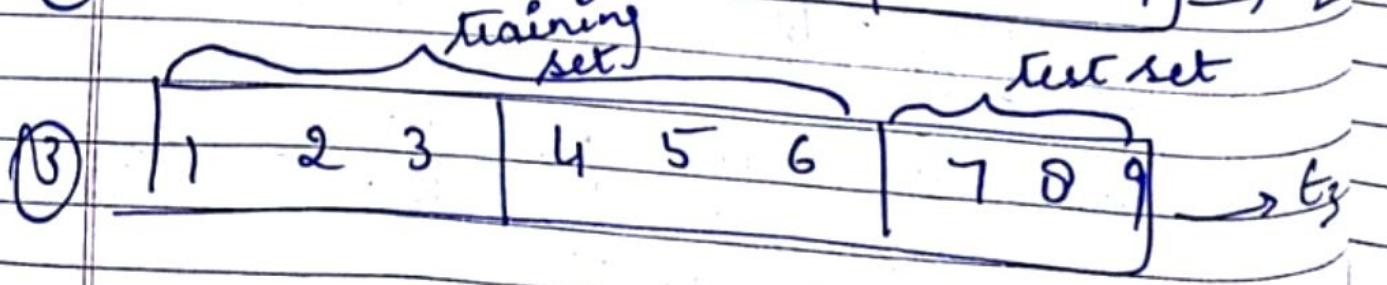
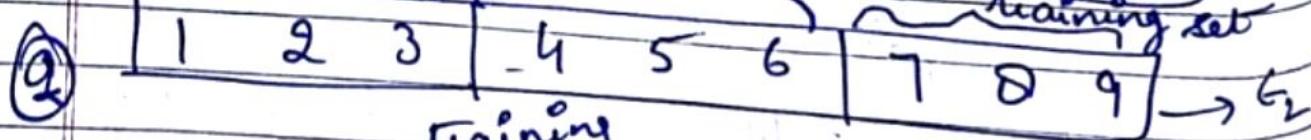
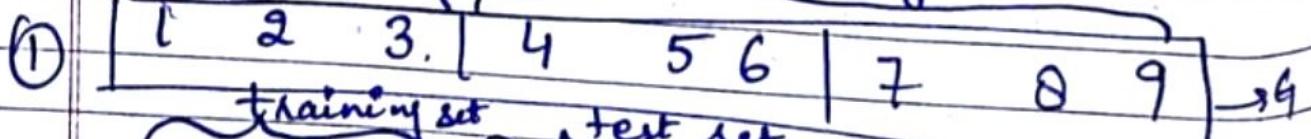
Cross-validation

systematic sampling

Not random sampling

→ It overcome the problem of "overlapping of test set".

divide the whole data into 3 subsets
3-fold cross-validation of equal size



$$E = \frac{1}{K} \sum_{i=1}^K E_i$$

Page

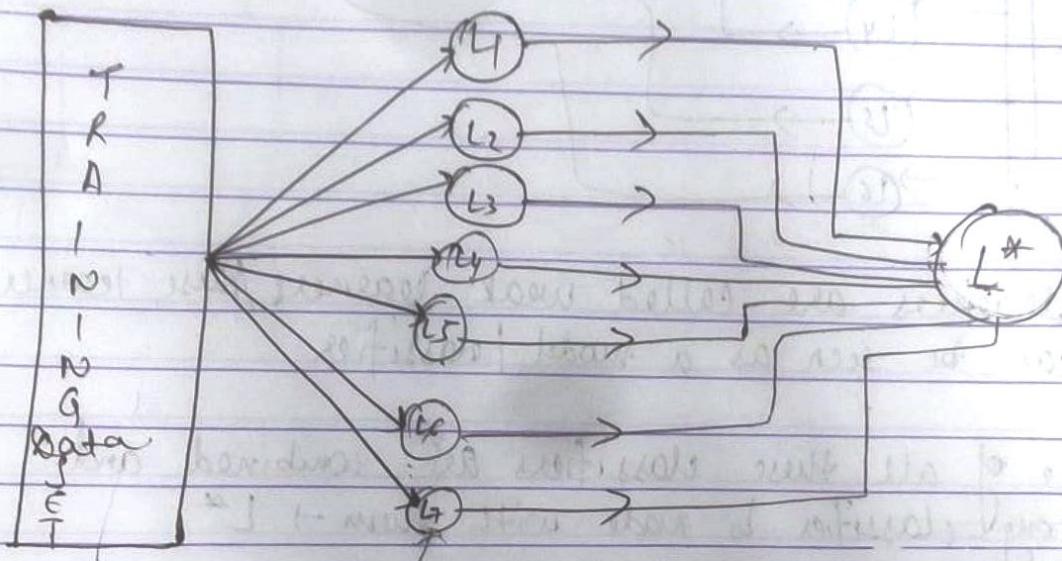
In k-fold cross validation, the initial data are randomly partitioned into K mutually exclusive subsets or folds " $D_1, D_2, D_3, \dots, D_K$ ", each of approximately equal size. Training set and testing is performed K times. In iteration i, partition D_i is reserved as the test set and the remaining partitions are collectively used to train the model.

- In the first iteration, subsets D_2, D_3, \dots, D_K collectively as the training set to obtain a first model, which is tested on D_1 .
- Second iteration is trained on subsets D_1, D_3, \dots, D_K and tested on D_2 and so on.
- Leave One-out is a special case of k-fold cross validation where K is set to the no. of initial tuples. That is, only one sample is "left out" at a time for the test set.

Ensemble Learning

Ensemble Learning: Don't take decisions based on only one single model.

→ Consider the output of different other models to reach to a decision/conclusion.



→ $L_1, L_2, L_3, L_4, \dots$ are the base learners.

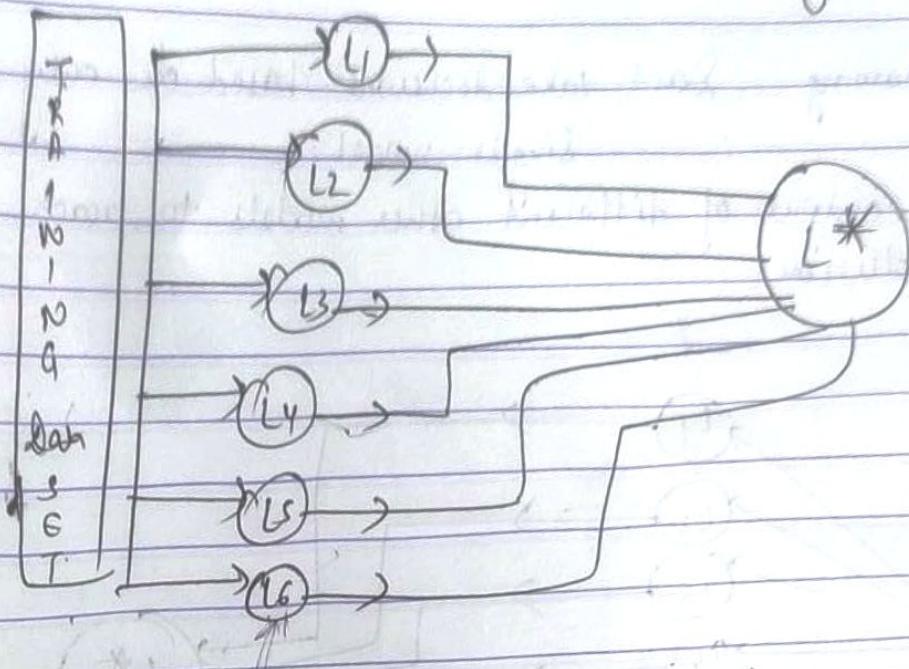
→ These learners can use different algorithms to maintain the diversity of their output.

→ This situation is called Heterogeneous situation.

→ If all the learners use the same algorithm, then the diversity of output will reduce.

→ To bring diversity in this case, we use different training data sets.

When learners use different Training data sets,



→ These learners are called weak learners / base learner.
+ This can be seen as a model / classifier.

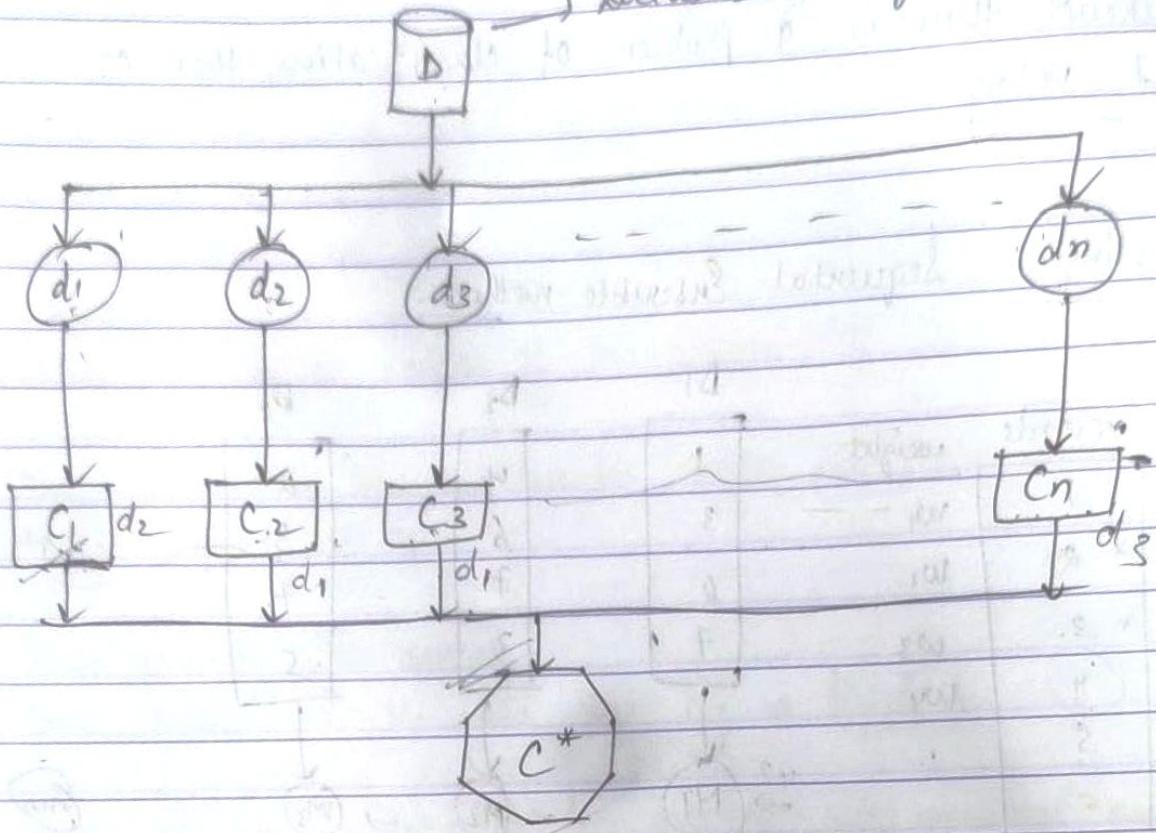
→ The o/p of all these classifiers are combined and a strong classifier is made with them → L^* .
→ Strong because the predictability of L^* as compared to other (L_1, L_2, \dots) is very high.

→ Three types of Ensemble Methods:

- ① Bagging
- ② Boosting
- ③ Random Forest

→ Bagging: Also known as Bootstrap Aggregation.

Database (Original Training Data).



→ $d_1, d_2, d_3 \dots d_n \rightarrow$ we randomly select samples from training data and create a new dataset with name d_i .

Similarly, we create new data sets $d_2, d_3 \dots d_n$.

This sampling is done with replacement. A record can be present in more than 1 datasets.

This dataset is also known as Bootstrap sample.

→ Now, classifiers are trained using these data sets.

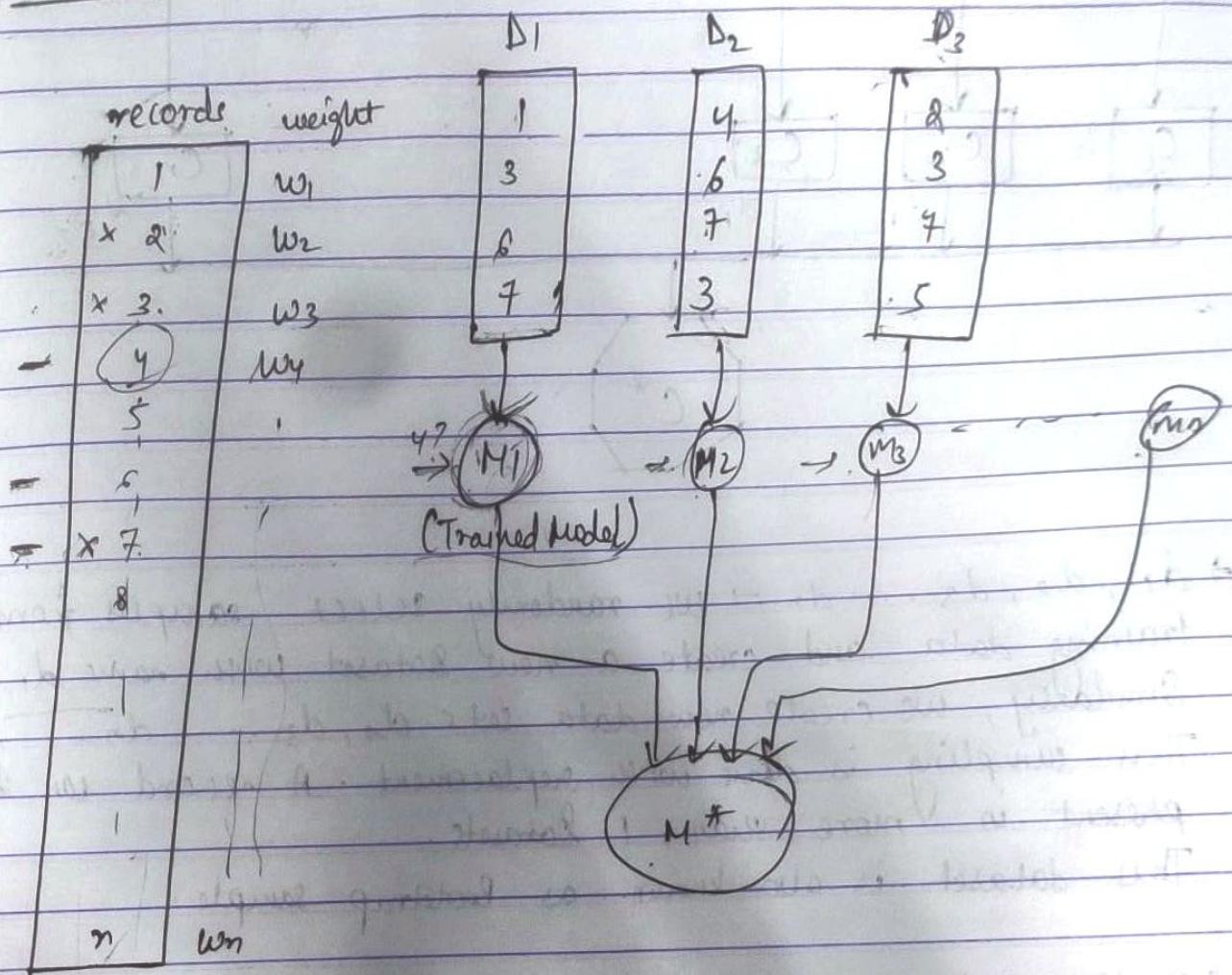
C_1 is trained from d_1 , C_2 is trained from d_2 , ... etc

Now, these classifiers are combined and a final ensemble classifier C^* is formed. (With low error rate & more accurate result).

If we need to find to which class a data belongs, all the classifiers are asked. The result is the majority votes that have come.

Whenever there is a problem of classification, then we use Voting.

Boosting : Sequential Ensemble Method.



We assume that all the weights are equal.

When we are creating a new dataset D_1 , then all the records are equally probable to be a part of D_1 .

M_1 will tell to which class the records belong. All the records will be checked one by one.

It is not necessary that M_1 also tells the correct data, to overcome this, 2nd model is trained.

Suppose $w_1, w_2, w_3 \rightarrow$ misclassified.

To move from one model to another, the weights need to be updated.

Focus more on the weights of the records that M_1 misclassified.

The ~~mis~~ records with misclassified weight now have more probability to get selected into D_2 .

Same process will repeat.

Same repeats with M_2 and then D_3 is filled.

Once all the models are formed, they are combined and a strong model M^* is formed.

In case of classification, Voting is preferred.

Random forest: uses decision trees.

It is a kind of ensemble classifier that uses decision tree algorithm

Initially, we must have original training data set.

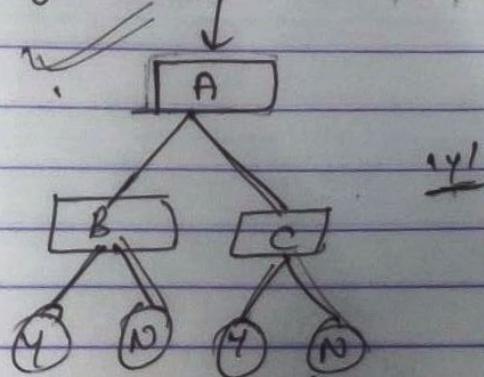
	A	B	C	T		A	B	C	T
1.				Y		2.			N
2.				N		3.			Y
3.				Y	→	4.			Y
4.				Y		5.			Y
5.				N					

Bootstrap Data Set

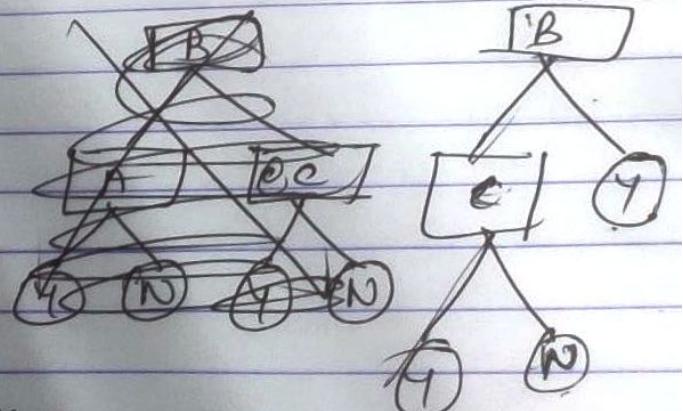
→ Create a Bootstrap data set using sampling.

Plot a decision tree using this Bootstrap Data Set

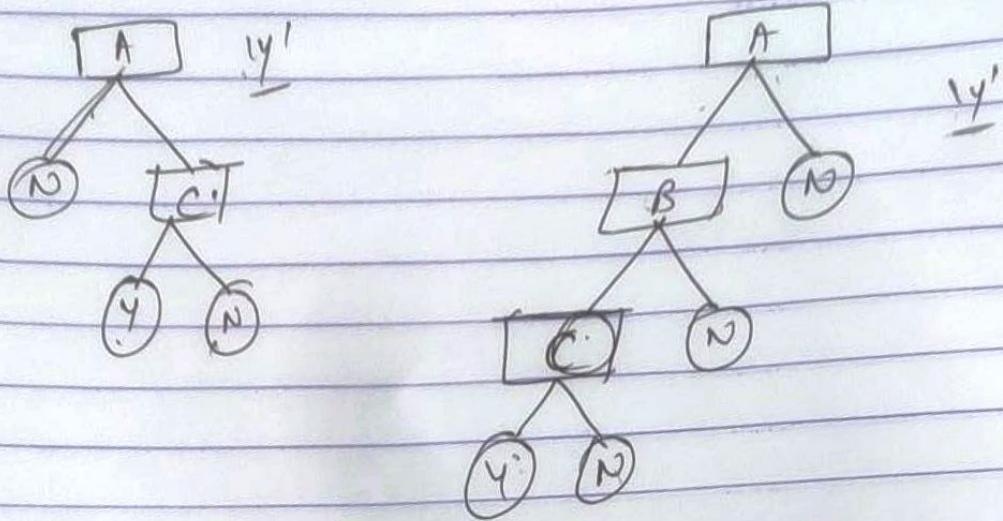
To decide the root node, out of 3 variables A, B, C we choose a subset of (A, B) and then assume that A is better for splitting.



same steps for (B, C)



Randomly generated Decision Trees.



Test tuple \rightarrow $[A | B | C | ?]$ \xrightarrow{T} $[A | B | C | ?]$ \xrightarrow{T}

find whether the target attribute will be Yes or No?

we apply this tuple to every decision tree.

we go with the ^{maxm} votes. $\xrightarrow{3Y}$
~~1N~~

Bootstrap Approach \rightarrow The training records are sampled with replacement

\rightarrow ensures accuracy.

OD

des	x	y	cls.
10	0	1	+
14	+	0	+
12	1	1	-

Sample

α_1	des	x	y	cls.
12	1	1	-	
12	1	1	-	
10	0	1	+	

similarity $\rightarrow \alpha_2, \alpha_3$

if original size is n , then sample will contain 63.21% of n

Bootstrapping Method - How does it work?

Invented by Bradley Efron, the bootstrapping method is known to generate new samples or resamples out of the already existing samples in order to measure the accuracy of a sample statistic.

Using the replacement technique, the method creates new hypothetical samples that help in the testing of an estimated value.

Here are 3 quick steps that are involved in the Bootstrapping method -

1. Randomly choose a sample size.
2. Pick an observation from the training dataset in random order.
3. Combine this observation with the sample chosen earlier.

For the samples that are chosen in the representative sample size, they are referred to as the ‘Bootstrapped samples’ or the bootstrap sample size. On the other hand, the samples that are not chosen are referred to as the ‘Out-of-the-bag’ samples that serve as the testing dataset.