



# Naive Bayes and Ensemble Technique HandWritten Notes

- Darshan R M

"Learn, Share,  
and  
Collaborate"

Week 17

Noise Bayes  
and  
Ensemble Technique  
and types

\* Noise Bayes algorithm  
    ↳ classification.

- ① Probability
- ② Bayes Theorem.

\* Events

    |  
    ③ Independent Events

    ↳ Probability of happening of  
    one event don't affect others.  
    ↳互不影响

Exe Rolling a dice  $\rightarrow \{1, 2, 3, 4, 5, 6\}$

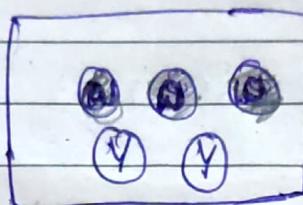
then

$$P(1) = \frac{1}{6}, P(2) = \frac{1}{6}, \dots$$

## (ii) Dependent events

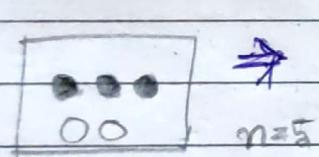
→ Probability of happening an event affect the happening of other event.

Ex:



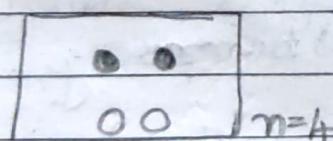
3 → orange  
2 → yellow

① what is the probability of removing a orange marble and then a yellow marble?



$$P(O) = \frac{3}{5} \rightarrow 1^{\text{st}} \text{ event}$$

Total orange marble  
Total marble.



$$P(Y) = \frac{2}{4} \rightarrow 2^{\text{nd}} \text{ event}$$

~~P(Y/O)~~

orange has been removed

Probability of first orange marble then yellow marble

$$P(O \text{ and } Y) = P(O) \times P(Y/O)$$

Conditional Probability:  $P(Y/O)$  = Probability of yellow marble when orange marble has been taken out.

$$P(\text{then } y) = P(\text{C}) * P(y|\text{C}) \\ = \frac{3}{5} \times \frac{1}{2} = \underline{\underline{\frac{3}{10}}}$$

~~generally~~

~~1<sup>st</sup> event~~      ~~2<sup>nd</sup> event~~

$$P(A \text{ and } B) = P(A) * P(B|A)$$

$P(A \text{ and } B)$  ~~Probability of~~ Probability of happening of event A and then happening of event B.

$P(B|A) \Rightarrow$  Probability of event B when event A has (completed) happened.

### \* ~~Baye's~~ Baye's theorem

w.k.t  $P(A \text{ and } B) = P(B \text{ and } A)$

$$P(A) * P(B|A) = P(B) * P(A|B)$$

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

$P(A|B) \Rightarrow$  Probability of event A when event B has already happened.

$P(A) \Rightarrow$  Probability of event A

\* Dataset independent feature      dependent feature

$x_1$	$x_2$	$x_3$	$y$
-	-	-	Yes
-	-	-	No
-	?	-	Yes
-	?	-	No

Now, we solve this entire prediction by using Bayes theorem.

~~The~~ ~~algorithm~~ which uses bayes theorem is called "Naive Bayes algorithm".

$$\Rightarrow P(y/x_1, x_2, x_3) = \frac{P(y) * P(x_1/x_2, x_3)/y}{P(x_1, x_2, x_3)}$$

Probability of  $y$  given that  $x_1, x_2$  and  $x_3$

$$P(y/x_1, x_2, x_3) = \frac{P(y) * P(x_1/y) * P(x_2/y) * P(x_3/y)}{P(x_1) * P(x_2) * P(x_3)}$$

feel "Yes"

\*  $P_{\text{Yes}} \left[ \begin{matrix} \text{Yes} \\ (x_1, x_2, x_3) \end{matrix} \right] = \frac{P(\text{Yes}) * P(x_1 \text{ yes}) * P(x_2 \text{ yes}) * P(x_3 \text{ yes})}{P(x_1) * P(x_2) * P(x_3)}$

constant

if  $\left[ \begin{matrix} y \\ 0.60, 0.40 \end{matrix} \right]$  then conclusion /  
output is "yes"

feel "No"

\*  $P_{\text{No}} \left[ \begin{matrix} \text{No} \\ (x_1, x_2, x_3) \end{matrix} \right] = \frac{P(\text{No}) * P(x_1 \text{ no}) * P(x_2 \text{ no}) * P(x_3 \text{ no})}{P(x_1) * P(x_2) * P(x_3)}$

constant,

so, we can remove  
this denominator.

a) Consider dataset

Day	Outlook	Temperature	Humidity	Wind	Rain	Temper
1	Sunny	High	High	Strong	No	Yes
2	Sunny	High	High	Weak	No	No
3	Overcast	High	High	Weak	Yes	No
4	Overcast	High	High	Strong	Yes	Yes
5	Rainy	Low	Low	Weak	No	No
6	Rainy	Low	Low	Strong	Yes	No
7	Rainy	High	High	Weak	No	No
8	Sunny	Low	Low	Strong	No	Yes

From dataset

	outlook		$P(E/\text{Yes})$	$P(E/\text{No})$
	Yes	No		
Sunny	2	3	$2/5$	$3/5$
overcast	4	0	$4/4$	<del><math>0/4</math></del>
Rain.	3	2	$3/5$	$2/5$
	$\Sigma = 9$	$\Sigma = 5$	$\Sigma = 1$	$\Sigma = 1$

### Temperature

	Yes	No	$P(E/\text{Yes})$	$P(E/\text{No})$
Hot	2	2	$2/4$	$2/5$
Mild	4	2	$4/6$	$2/5$
Cool	3	1	$3/4$	$1/2$
	$\Sigma = 9$	$\Sigma = 5$		

tell probability  
of Cool temperature  
when play tennis  
per No.

	<u>Play</u>	$P(\text{Yes})$	$P(\text{No})$
Yes	9	<del><math>9/14</math></del>	<del><math>5/14</math></del>
No	5	$9/14$	$5/14$

if:  $\frac{?}{8}$        $\frac{?}{8}$   
 Test (Sunny, Hot)  $\rightarrow$  ?

$$\Rightarrow P\left(\frac{\text{Yes}}{\text{(Sunny, hot)}}\right) = \frac{P(\text{Yes}) \times P(\frac{\text{Sunny}}{\text{Yes}}) \times P(\frac{\text{Hot}}{\text{Sunny}})}{P(\text{Sunny}) \times P(\text{Hot})}$$

$$= \frac{\frac{3}{14} \times \frac{3}{9} \times \frac{2}{5}}{\frac{2}{63}} = \underline{\underline{0.031}}$$

$$\Rightarrow P\left(\frac{\text{No}}{\text{(Sunny, hot)}}\right) = \frac{P(\text{No}) \times P(\frac{\text{Sunny}}{\text{No}}) \times P(\frac{\text{Hot}}{\text{Sunny}})}{\text{constant}}$$

$$= \frac{\frac{1}{14} \times \frac{3}{9} \times \frac{2}{5}}{\frac{3}{63}} = \underline{\underline{0.085}}$$

$$\Rightarrow P\left(\frac{\text{Yes}}{\text{(Sunny, hot)}}\right) = \frac{0.031}{(0.031 + 0.085)}$$

$$= 0.27$$

$$= \underline{\underline{27\%}}$$

$$P\left(\frac{\text{No}}{\text{(Sunny, hot)}}\right) = \frac{0.085}{0.031 + 0.085}$$

$$= 0.73$$

$$= \underline{\underline{73\%}}$$

<u>outlook</u>	<u>Temperature</u>	<u>C/F</u>	
Sunny	Hot	73%	Not play Tennis
		27%	They will <del>not</del> play Tennis

↓

① person is not playing Tennis

### \* Variants of Naive Bayes

- ① Bernoulli Naive Bayes.
- ② Multinomial Naive Bayes.
- ③ Gaussian Naive Bayes.

### ① Bernoulli Naive Bayes :

⇒ whenever your features are following a Bernoulli distribution, then we need to use Bernoulli naive bayes algorithm.

When there is only 2 outcome [0/1]

## Dataset

$f_1$	$f_2$	$f_3$	$f_4$	$O/P$	$f_1$	$f_2$	$f_3$
Yes	Pass	Male	Yes	1	1	1	1
Yes	Fail	Female	No	1	0	0	0
No	Pass	Male	Yes	0	1	1	1
Yes	Fail	Female	No	1	0	0	0

all features following  
Bernoulli distribution

↓  
Sparse Matrix

weak

use NLP technique

## ② Multinomial Naive Bayes

$\Rightarrow$   $O/P \Rightarrow$  raw form of Text  
and

$O/P \Rightarrow$  classification

## Dataset's Spam classification.

Email	Spam / Not spam
Baby	
You won 100 Lakh You are good	Spam Ham.

Convert into "numerical"  
values

using  
NLP  
[Natural]

Inside we will  
have a formula of  
which converts into  
numerical.

① Bow

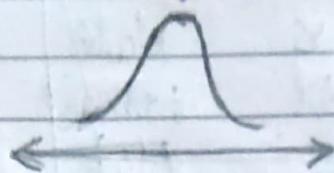
② Tfidf

③ Word zvec.

← language  
processing]

### ③ Gaussian Naive Bayes

If the features are following Gaussian distribution, then we use gaussian naive bayes.



Iris dataset

features are continuous

Age Height weight

	Age	Height	weight	Overweight
	25	170	78	Yes
	38	160	75	No
	40	130	60	Yes
	51	170	35	No
	{	}	{}	{}

\* we can use Bernoulli and multinomial naive bayes in NLP techniques

## \* Implementation

- > from sklearn.datasets import load\_iris
- > from sklearn.model\_selection import train\_test\_split
- >  $x, y = \text{load\_iris}(\text{return\_X\_y} = \text{True})$

① > train-test split  $\rightarrow$  30%

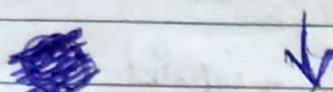
② Model Building

> from sklearn.naive-bayes import GaussianNB

- > gnb = GaussianNB()
- > gnb.fit(x\_train, y\_train)
- > gnb.predict(x\_test)

③ Performance metrics

\* Azure deployment  
    Create Web App Service.



Configuration required for  
github will be taken.  
Care by azure.

Can see

\* github → actions ↴

[Build and deploy  
section]

\* Ensemble Techniques and Bagging

↳ Used to increase accuracy  
and efficiency of the model

① What is Ensemble ?

⇒ Combining multiple → Train → Prediction  
models

Two types

① Bagging.

Technique

(i) Random forest classifier and regressor

② Boosting.

Technique

(i) Ada Boost.

(ii) Gradient Boost.

(iii) Xg Boost.

Bagging Technique

→ used in hackathons  
without much hyperparameter tuning model performs well.

Dataset

Sample data

$d''$

M1

Decision Tree  
classifier

D

$d'''$

M2

Naive Bayes

$d''''$

M3

logistic regression

"m" rows

$d'''''$

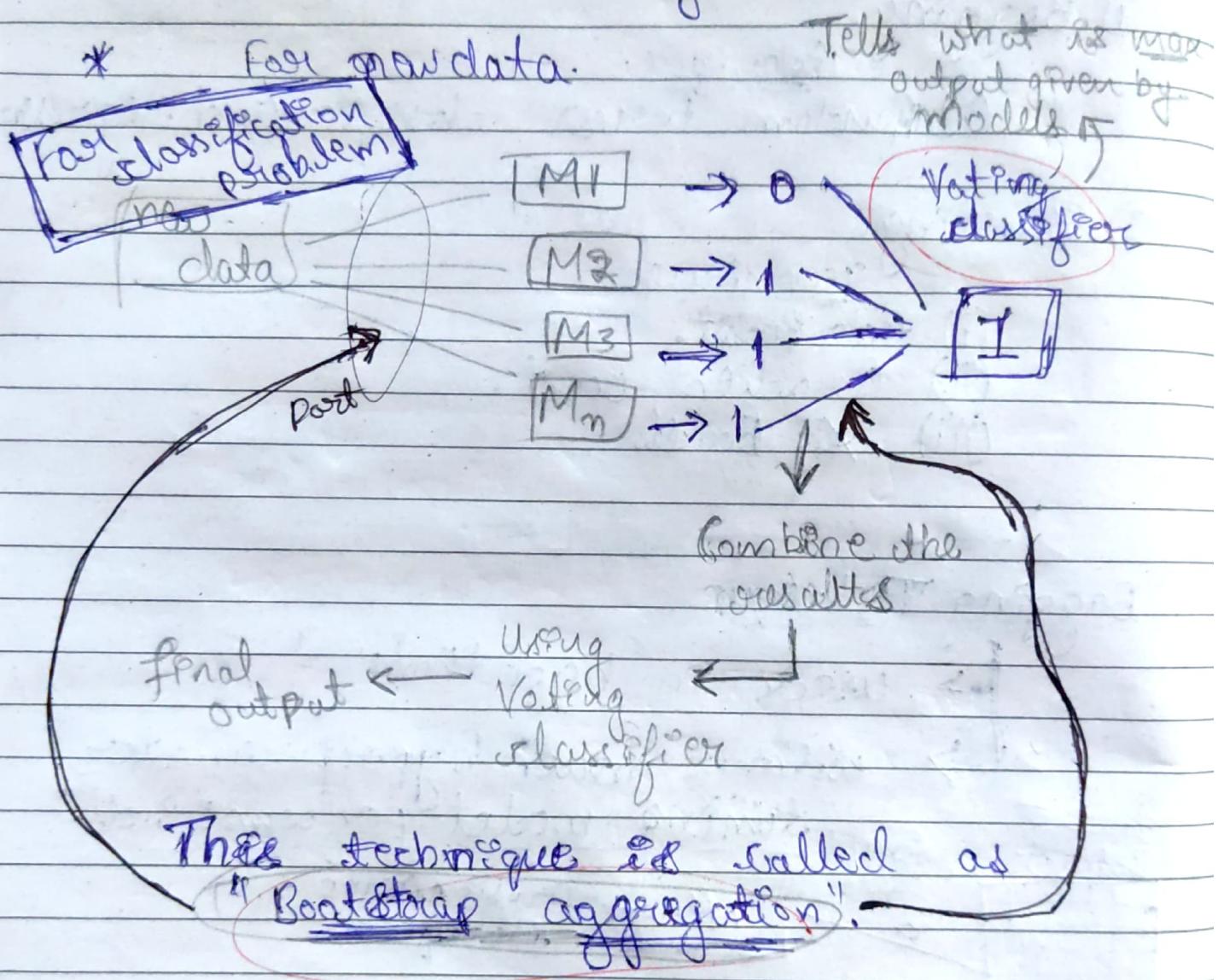
M $m$

Naive Bayes

"m" independent  
models

Basis learners

- \* After training ~~the~~ basic learner with those corresponding sample data.



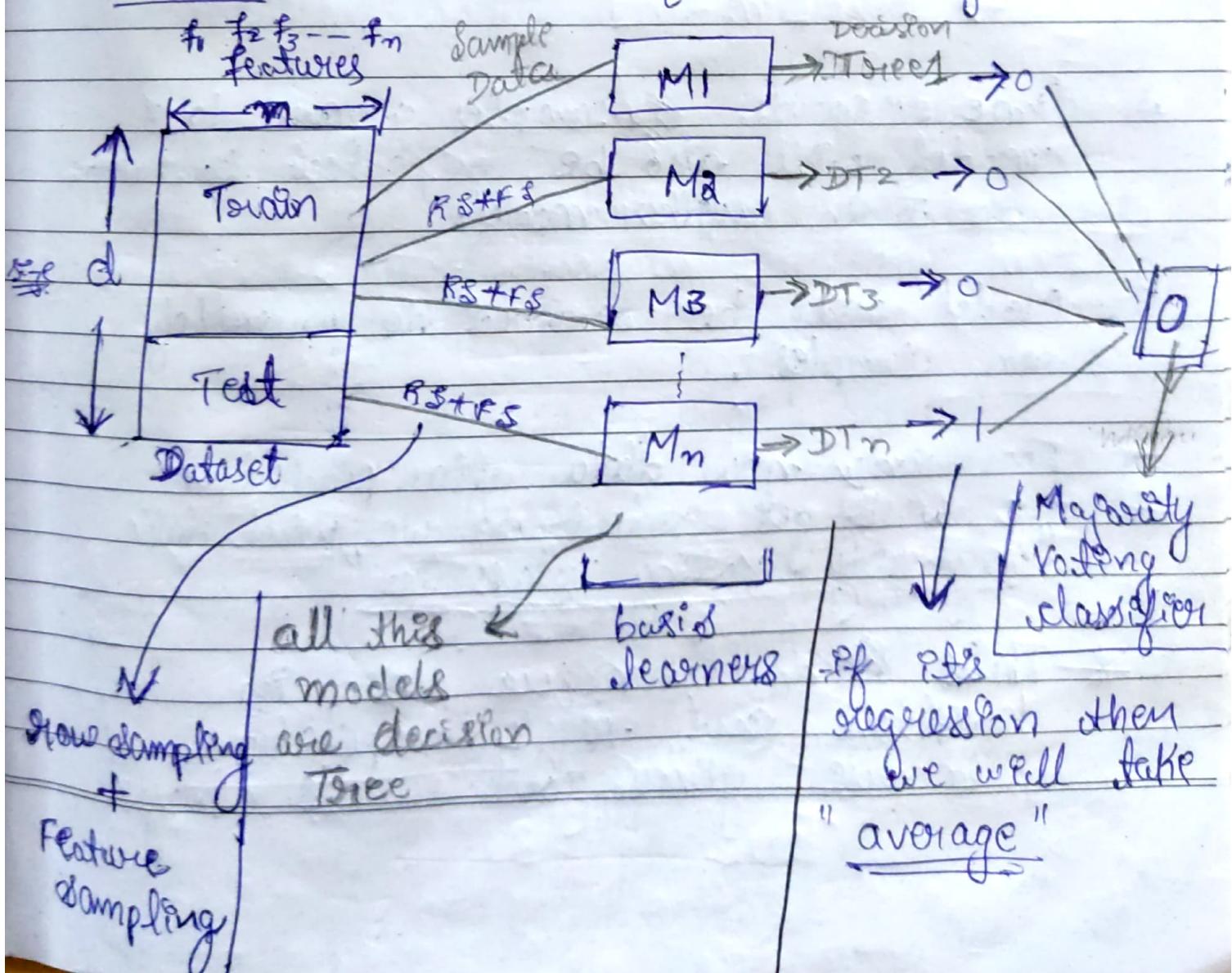
- \* ~~For~~ all the base learners of random forest we use decision Tree.
- \* The technique where we combine different different model is called "Custom bagging technique"

\* For regression problem.

basis learner  $\xrightarrow{\text{will be}}$  regression model

we get output for continuous data  
then we will take average  
of all ~~all~~ output to do the prediction.

## \* Random Forest classifier and regressor



\* Row-sampling  $\rightarrow$  from the train data we take some records and give to model for training.

if  $d'$  samples  
 $\hookrightarrow$  then  $[d' < d]$

\* Feature sampling  $\rightarrow$  from the train data we take particular feature for model training.

if  $m'$  samples  
 $\hookrightarrow$  then  $[m' < m]$

\* Some records. ~~samples~~ given to 1 model will also be repeated to taken to train another model.

\* Model will be trained to parallelly on samples.

\* For any new data after prediction we will get output so we will take "majority voting classifier"

\* This base technique used in Hackathon and Kaggle competition. because gives good results always.

Notes:

Classification  $\rightarrow$  majority voting classifier.

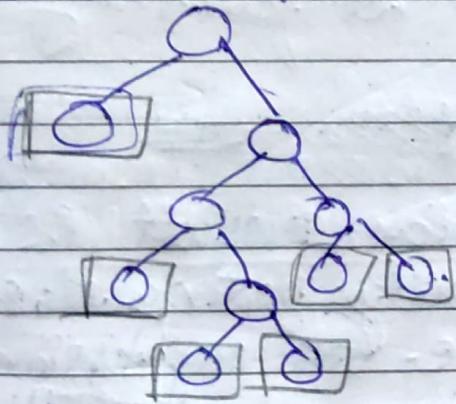
Regression  $\rightarrow$  average - o/p of all the models.

VIMP

Q) why should we use random forest instead of decision tree?

$\Rightarrow$  In DT we can also do post pruning and pre-pruning [using hyperparameter tuning]

Decision Tree



$\Rightarrow$  we take feature and split until we get leaf node.

but if we split until leaf node it leads to "overfitting".

Training  $\rightarrow$  acc<sup>↑↑</sup> low bias and test  $\rightarrow$  acc<sup>↓↓</sup> high variance

but To get "generalized model", we need to get "low bias and low variance".

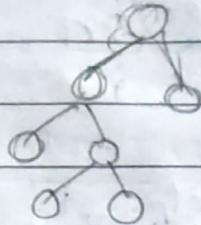
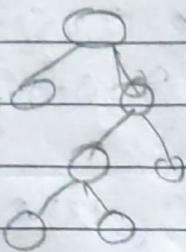
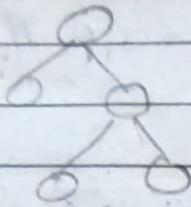
Possible only using "random forest"

different decision  
Tree

\* In Random forest Basic learners are trained parallelly.

Generalized model

↳ we get parallel expertise



These are different decision Tree which are good at few thing

~~all while predicting~~

do our prediction result is taken from ~~will be~~ majority voting classifier / ~~and~~ average (In regression) we get good result

[IMP]

\* Out of Bag Score Decision Tree ↴

interview  
asked

In Random forest ~~when~~ while training model we may get a ~~some~~ of missing some records for training that data is called "Out of Bag Data". To use that data we have a parameter ~~oob score~~ when it is True Out of Bag data will be used ~~at~~ Validation data to calculate accuracy and performance of model

Dataset

Train

Test

out of  
Bag  
data

that  
data  
is called

Some part of training data  
will not be selected for  
training.

why?

As we take random  
sample from data.  
there is a possibility  
that some record  
will not be selected  
for training.

if this data is not select for training  
but is learned. because of random  
sampling / feature sampling

we are losing this  
in term ← data.

we are hammering the  
performance of model by missing  
this data.

\* So, while performing random  
forest will get this scenario.

\* There is 1 parameter  
out of bag score  
oob score → Tree



then out of Bag data will be  
considered as "validation data"  
and also used to check  
"accuracy & performance"



These results are  
called as "oob-score"

Can be: 85%, 90%, 75%

In python,

> from sklearn.ensemble import  
RandomForestRegressor/RandomForestClassifier  
parameters

① n\_estimators → No of trees in forest  
→ 100 default

② criterion → entropy, gini, log. loss

③ oob\_score. → (Default → false)

|  
many others

|  
, out of Bag data  
ignored  
by default.

~~Production grade application~~. understanding.

~~Q&A~~

## \* Implementation

with

Pipelines



Vamp:

whole building production  
grade application.

and Hyperparameter  
tuning

Until, we were doing Feature selection,  
engineering and handling different value  
then training the model.

this model training is not one  
time activity because when  
we get "new data" we will  
refresh the model.

so, we need  
to write  
code when

can automate  
the entire process

→ so, for that we cannot  
keep on do feature engineering  
separately.

① TPS dataset.

② Label Encode → 'Time' feature

③ Segregate ~~Feature Test~~ independent  
/ dependent,

④ Train Test split.

missing values  
→ from sklearn.impute → import SimpleImputer

→ from sklearn.preprocessing → OneHotEncoder  
Categorical → StandardScaler  
Numerical → Feature Scaling

## ⑤ Pipeline

\* we can automate process.

\* we can combine multiple steps

> from sklearn.pipeline import Pipeline

> from sklearn.compose import

ColumnTransformer

used to combine several feature extraction mechanism into single transformer.

(i) Identify categorical and numerical columns.

> cat-cols = [ ]  
> num-cols = [ ]

(ii) feature engineering automation.

> main\_pipeline = Pipeline(

steps = [

'imputer', SimpleImputer(strategy='median'))  
'scaler', StandardScaler()])

]

name

↑

object

to be done

~~>~~ Cat-pipeline = Pipeline ( steps = [ "most frequent"  
 { "imputer", SimpleImputer (strategy = ↓ ) },  
 { "onehotencoder", OneHotEncoder (P) },  
 ] )

~~&~~ Now, as we use both pipeline we have to ~~wrap~~ them with ColumnTransformer

~~wrapper~~  
~~>~~ preprocessor = ColumnTransformer ( [ name pipeline name features name ]

( "num-pipe", num_pipeline, numerical cols ),		
( "Cat-pipe", CatPipeline, categorical cols )		

~~x-train~~ = preprocessor.fit\_transform  
 (x\_train)

## ⑥ Model Building + Evaluation

~~>~~ from sklearn.ensemble import  
 RandomForestClassifier.

~~&~~ Automating model training process.

\* models = {

} 'Random Forest': RandomForestClassifier(),

\* function which calculate the do model scoring and evaluation

~~def~~ def evaluate\_models(x-train, y-train, x-test, y-test, models):

report = {}

for i in range(len(models)):

Taking model. model = list(models).~~values()~~[i]

training model. fit (x-train, y-train)

predicting y-pred = model.predict(x-test)

Calculating accuracy model-score = accuracy\_score(y-test, y-pred)

storing accuracy. report[~~list~~(models.keys())[i]] = model-score

return report

## ④ Hyperparameter Tuning

> from sklearn.model\_selection import  
RandomizedSearchCV

> Rand\_cv = RandomizedSearchCV(  
RandomForestClassifier(),  
param\_distributions = parameters,  
cv = 3, scoring = "accuracy",  
verbose = 3)

> Rand\_cv.fit(x\_train, y\_train)

Internal assignment

off : total bill.

## \* AdaBoost

Ensemble  
Technique

Decision  
tree.

Bagging

{ Basis learner }

- ① Random Forest classifier.
- ② Random Forest Regressor.

Boosting

{ weak learner }

- ① Ada Boost.
- ② Gradient Boost.
- ③ XG Boost.

{ basis learner } → Decision  
Trees

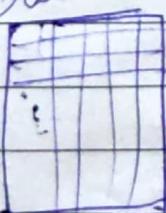
low Bias &  
low variance

Can solve both  
classification and  
regression  
problems.

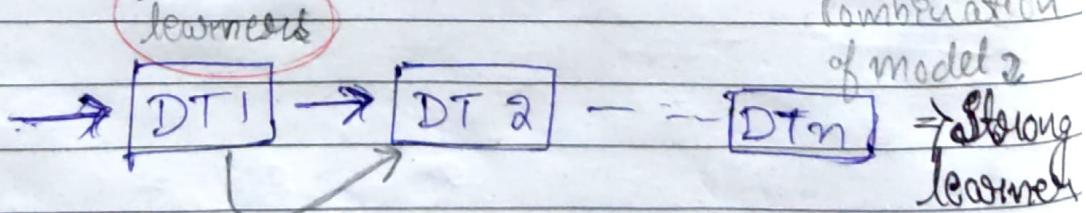
## \* Boosting

{ All models are sequentially connected }

Dataset



weak  
learners



Combination  
of model 2

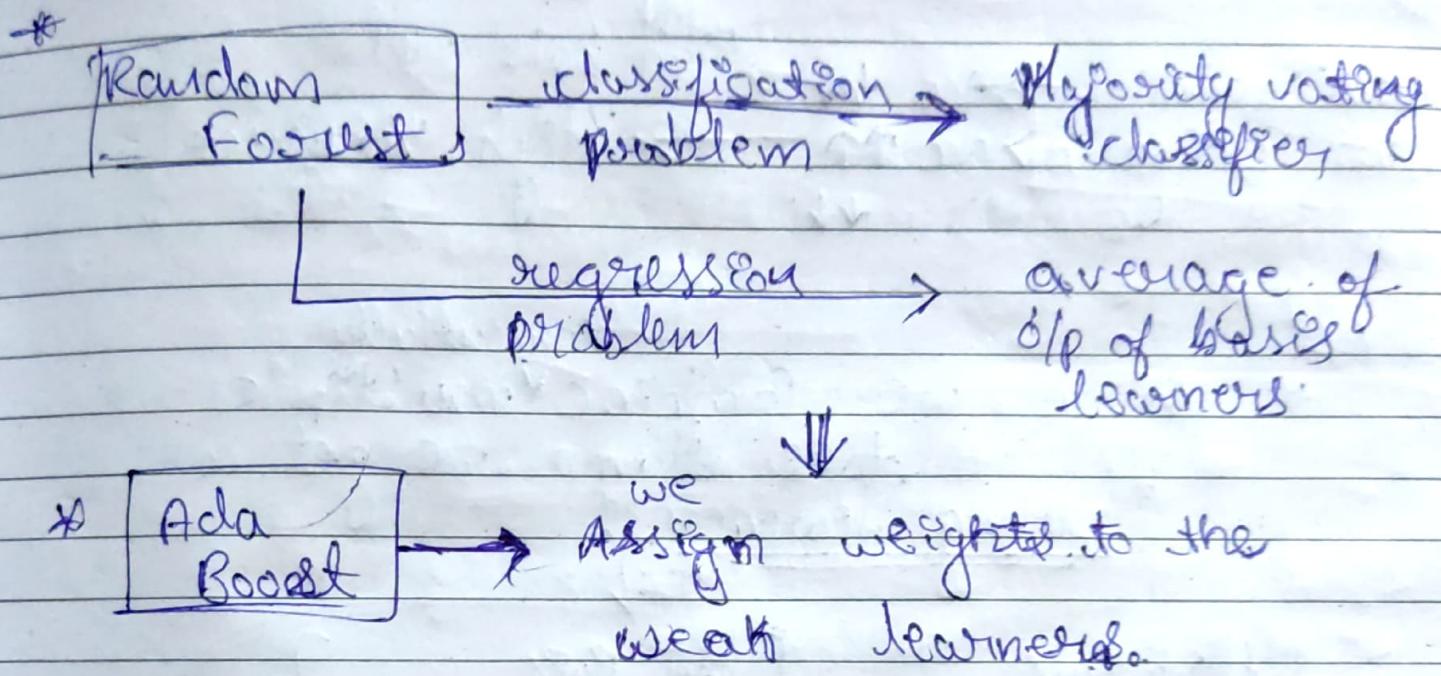
Strong  
learner

wrongly  
predicted

output send  
to next decision

Trees.

\* weak learner  $\rightarrow$  haven't learnt much from the training dataset.

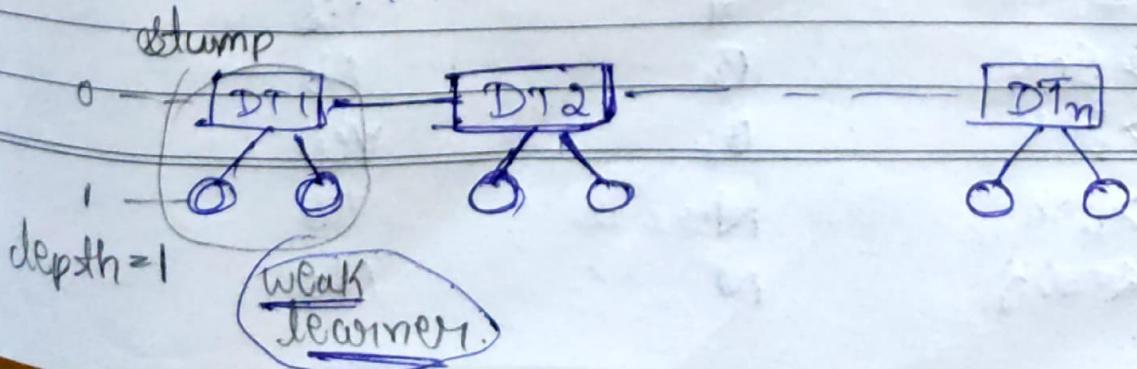


$$f = d_1(M_1) + d_2(M_2) + \dots + d_n(M_n)$$

both  $\{M_1, M_2, M_n\}$   $\rightarrow$  Decision tree stumps ~~stumps~~  
 classification and regression  $\{d_1 \text{ to } d_n\} \rightarrow$  weights.

\* Decision Tree stump

a A tree with 1 depth.



\* weak learners [Decision tree stump]

→ Under fitting ↴

40% Training → acc ↓ → High bias and  
45% Testing → acc ↑↑ → High Variance  
(or) acc ↓↓

\* when we combine multiple decision tree stump we get a strong learner.

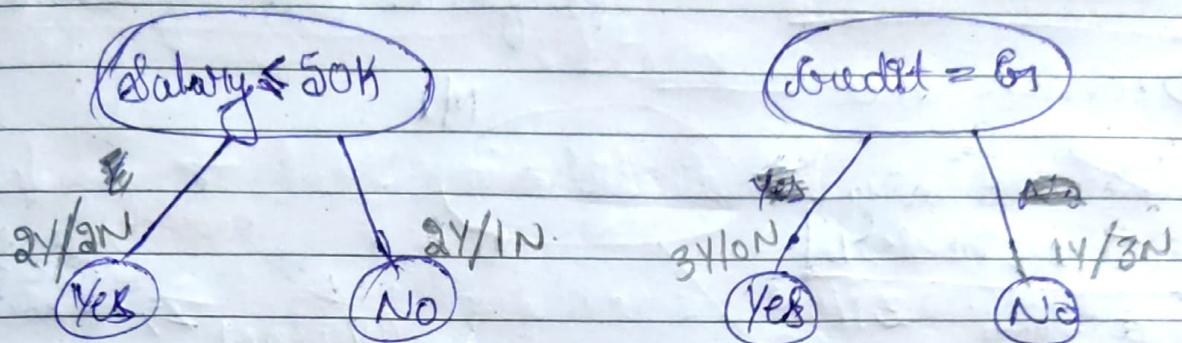
{ High Bias }  
High variance } → { low Bias }  
Decision Tree stump      High variance }

\* Ada Boost math

Salary      Credit & Care      Approval

$\leq 50K$	Bad	No
$\leq 50Y$	Good	Yes
$\leq 50B$	G	Yes
$> 50K$	B	No
$> 50K$	G	Yes
$\leq 50K$	Normal	Yes
	N	No

~~Step 1~~ we create Decision Tree stump ~~&~~ select best one



DT stump-1

DT stump-2

- \* we can create other DT stumps also
- \* Now, From above 2 stumps which one we have to select ?

we use Impurity check

Entropy  $\leftarrow$   $\rightarrow$  Gini

$$H(S) = -(P_+)(\log P_+) - (P_-)(\log P_-)$$

$$Gini = 1 - \sum_{i=1}^n P_i^2$$

- \* In DT stump-1  $\rightarrow$  Yes  $\Rightarrow$  1 [Impure]  
No  $\Rightarrow$   $\approx 1$  [-1-]

In DT stump-2  $\rightarrow$  Yes  $\Rightarrow$  0 [pure split]  
No  $\Rightarrow$  0. - [less impure split]  
Tells pure leaf node easily  $\leftarrow$  good

- \* we select best decision tree stump using entropy

step 1 overview

Step 2: Sum of the total error and performance of stump.

Decision tree

which is selected using entropy in step 1

Credit is "good" 3Y/ON  
the option all Yes.

MI

Credit = G<sub>1</sub>

1 Yester

Error

Credit is not good but we have 1 Yes %.

Helps to calculate sum of total error

library condition Action sample weights

≤50K	R	No	1/7
≤50K	G <sub>1</sub>	Yes	1/7
≤50K	G <sub>1</sub>	Yes	1/7
>50K	R	No	1/7
>50K	G <sub>1</sub>	Yes	1/7
>50K	N	Yes	1/7
≤50K	N	No	1/7

Deviations

assign equal weights

only 1 wrong record

① Total error = 1/7

sample weights

$$\textcircled{2} \quad \text{Performance of Stump} = -\frac{1}{2} \log_e \left[ \frac{1 - TE}{TE} \right]$$

$$\begin{aligned} \text{Performance of Stump} &= -\frac{1}{2} \ln \left[ \frac{1 - TE}{TE} \right] \\ &= \frac{1}{2} \ln [6] \end{aligned}$$

$$\approx \underline{\underline{0.896}}$$

(indicates weight.)

~~if performance~~

what does it indicate?

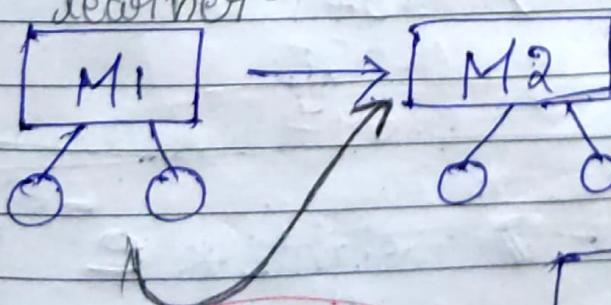
$d_1$

WKT

$$f = d_1(M_1) + d_2(M_2) + \dots + d_n(M_n)$$

$$d_1 = 0.869 \rightarrow \text{weight}$$

\* weak learner [Decision Tree Stump]



wrongly predicted records

overview of step 2

\* we calculate performance of stump using Total error. tells weight to assign to stump

\* aim: To send incorrigible misclassified points to next DT stump so we have in it's probability of selecting by updating weight

Step 3: update the weights for correctly and incorrectly classified points.

\* \* Correctly classified  $\rightarrow$  decrease selecting probability weight

we have to in correctly classified  $\rightarrow$  increase weight

because of this there is a high probability of selecting wrong record by next DT stump.

	<u>Simple weights</u>	<u>Updated weights</u>
	$\frac{1}{7}$	0.058
<del>incorrectly classified</del>	$\frac{1}{7}$	0.058
<del>incorrectly classified</del>	$\frac{1}{7}$	0.058
250	N [ Yes ] $\frac{1}{7}$	0.349
	$\frac{1}{7}$	0.058
	$\frac{1}{7}$	0.697

\* For correctly classified points  $\rightarrow$  weights  $\times e^{-\text{Performance of stump}}$

$$= \frac{1}{7} \times e^{-(0.596)}$$

$$= \underline{\underline{0.058}}$$

For incorrectly classified points  $\rightarrow$  weights \* e<sup>Performance of stump</sup>

$$= Y_1 \times e^{0.896}$$

$$= \underline{\underline{0.349}}$$

Step 4: Note  
Normalized weights computation  
and assigning bias.

Normalized weight

\* we can notice that updated weights sum is not equal to 1

observation  $\sum$  updated weights  $\neq 1 \leftarrow 0.697$   
step 3

but initially when we assign the weight

step 2 initially  $\rightarrow \sum$  sample weights  $\approx 1$

so, we have to make  $\sum$  updated = 1  
weight

so we use to normalize the weight

① updated weights

Normalized weights

weight  
eweights

Bonus  
assignment

0.058

%0.697

0.08

0.058  
0.697

0 - 0.08 ·

0.058

0.08

0.08 - 0.16

0.058

0.08

0.16 - 0.24

0.058

0.08

0.24 - 0.32

0.058

0.08

0.32 - 0.40

0.349

0.50

0.40 - 0.90

0.088

0.08

0.90 - 0.98

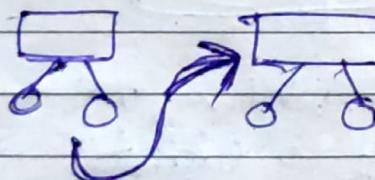
$$\Sigma = \underline{\underline{0.697}}$$

$$\Sigma = \underline{\underline{1}}$$

To make

Max  
Bin size

② bias assignment



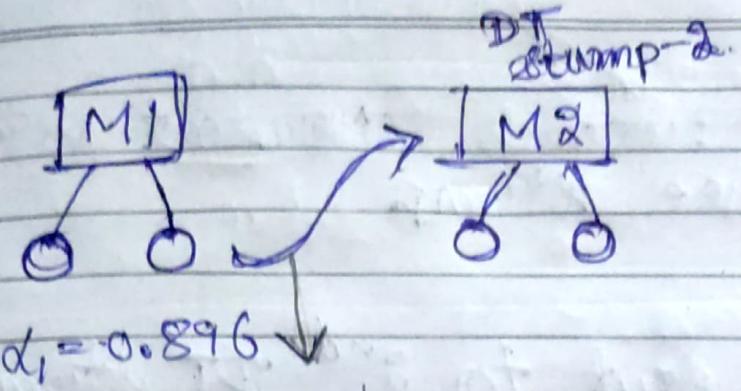
How to make our ml algorithm

only send "incorrectly  
classified points"?



and we "assign bins"

⇒ Now, we will send the incorrectly  
classified record frequently to next DT  
our (DT stump - 2)



we have prepare  
which datapoint/record  
do be sent

Step 5: Select data points to send to  
next stump

① Iterative process selecting random  
value between 0 and 1.

selected  
records

for next stump.

Bins	Random numbers	is	Credit Approval
assigns	[0 - 1] check bin and randomly generated	select record	
	0.50	$> 50K$	N Yes
(0.4 - 0.9) ↘	0.10	$\leq 50K$	G Yes
	0.60	$> 50K$	N Yes
	0.75	$> 50K$	N Yes
	0.24	$\leq 50K$	G Yes
	0.32	$\geq 50K$	B No.
	0.87	$> 50K$	N Yes

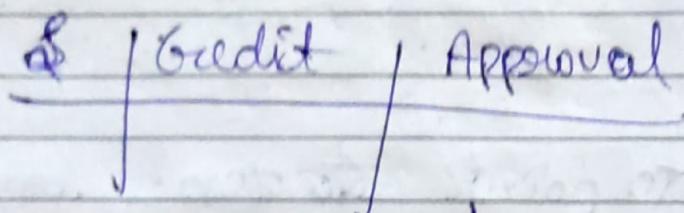
observation: incorrectly

classified records is  
picked again and again. as "ben size"  
is more as its weight is increased.

which contain more  
or correctly classified  
records.

- \* Now, "selected records" are sent to next decision tree stump.

Step 6: Record sent to  $2^{nd}$  DT stump.



again follow  
repeat same steps

→ If then

Performance stump  $\Rightarrow 0.65$

$$f_1 = d_1(M_1) + d_2(M_2)$$

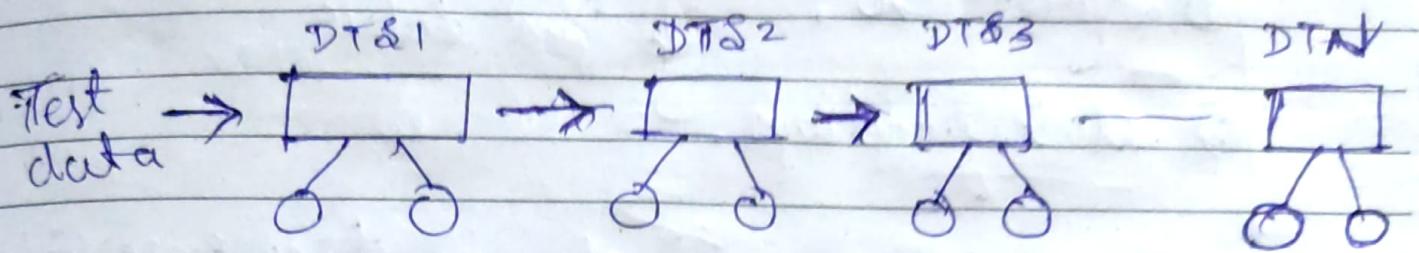
$$d_1 = 0.856 ; d_2 = 0.65$$

Step 7: Final prediction

Here, for classification problem

- \* If classification, we use 'mean square error' in place of 'entropy'  
and at final prediction instead of we get continuous value

Test - data. ( $\leq 50k$ ,  $G_1$ )



Boosting we will  $\downarrow$   
not use  $\rightarrow$  Yes  
voting classifier

~~Consider~~  $d_1 = 0.896$

performance of [denotes weight]  $d_2 = 0.65$

$d_3 = 0.24$  —  $d_n = -0.30$

$$f = d_1(M_1) + d_2(M_2) + \dots + d_n(M_n)$$

$$\Rightarrow (0.896)(\text{Yes}) + (0.65)(\text{No}) + (0.24)(\text{Yes}) \\ + (-0.30)(\text{No})$$

$$\Rightarrow (1.136)(\text{Yes}) + (0.35)(\text{No})$$

Performance of day (Yes) = 1.136

Yes  $>$  No

Performance of day (No) = 0.35

Final output  $\rightarrow$  Yes

## \* Gradient Boosting algorithm

- ① Regression.
- ② Classification.

→ It is a part of boosting ensemble technique. where in we create decision tree sequentially to get a "strong learner".

Consider, Regression dataset

Exp	Degree	Salary	$\hat{y}$	$y - \hat{y}$ Residual
2	B.E	50K	75K	-25K
3	Masters	70K	75K	-5K
5	Masters	80K	75K	+5K
6	PH.D	100K	75K	+25K
		$\mu = 75K$		

independent

dependent

[Continuous data  $\rightarrow$  regression problem]

Steps:

model should  
be unbiased

[and give  
default  
value]

① Create a Base-model

model

75K

average = 75K

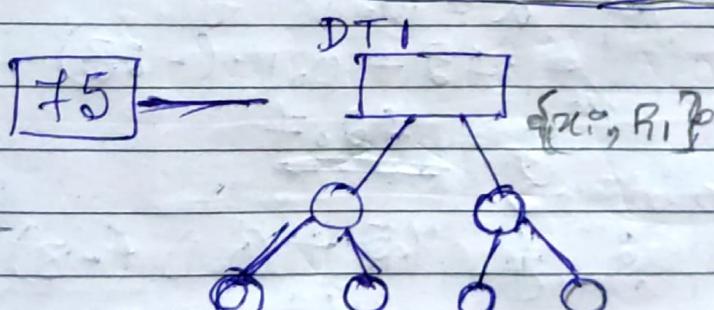
For any input we get 75K  
as output

② Compute residuals, Errors calculation

[actual - predicted]

③ Construct of Decision Tree

Consider s/p "x<sub>1</sub>" and  
o/p "R<sub>1</sub>"

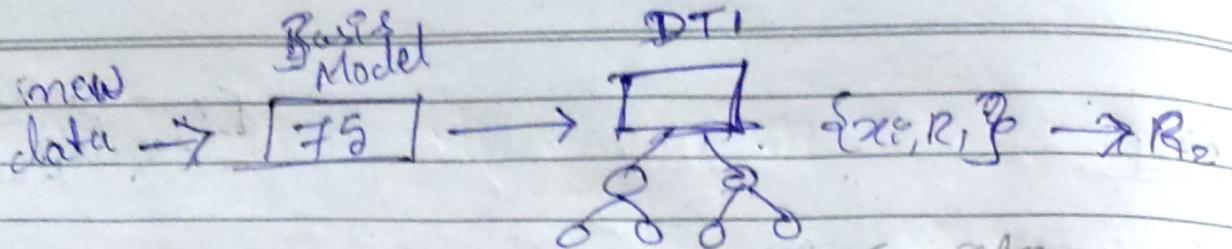


output of the  
decision tree  
of "R<sub>2</sub>"

"consider" R<sub>2</sub>

-23  
= 3  
+ 3  
+ 20

base +  
d(DT<sub>1</sub>)  
after we go  
and compute  
residuals e.g.  
to step 2



\* predicted o/p  $\Rightarrow$   $75 + (-23)$   
 $= \underline{\underline{52}}$

Is overfitting?  
 do not produce

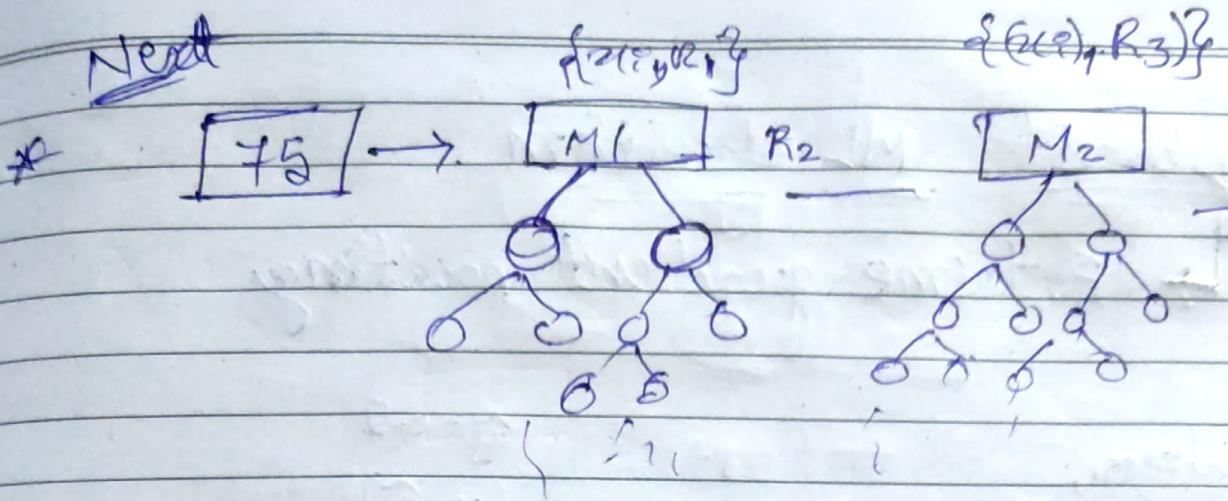
Model performing good at  
 actual value  $\rightarrow 50K$   
 predicted value  $\rightarrow 52K$

$\leftarrow$  Predicted o/p  $\Rightarrow 75 + \alpha \cdot (\text{DT}_1)$

$\alpha = \text{learning state } [0 \text{ to } 1]$   
 if  $\alpha = 0.1$   
 $= 75 + 0.1 \cdot (-23)$   
 $\stackrel{1^{\text{st}} \text{ term}}{=} \underline{\underline{72.7}}$

Calculation  
 $\stackrel{2^{\text{nd}} \text{ term}}{\rightarrow} \stackrel{\text{score}}{\rightarrow} 75 + (0.1) (-3)$   
 $= 75 - 0.3$   
 $= \underline{\underline{74.7}}$

<u>Y</u>	<u>R<sub>0</sub></u>
72.7	-22.7
74.7	-14.7
74.7	74.7



\* Step . Conclusions

\* Result

$$F(x) = h_0(x) + d_1 [h_1(x)] + d_2 [h_2(x)] + \dots + d_n [h_n(x)]$$

Learning state,  $d = 0.1$  [For all  $d_1, d_2, \dots, d_n$  we can apply]

$$\boxed{F(x) = \sum_{i=0}^n d_i [h_i(x)]} \Rightarrow \text{Final function}$$

\* we can do prepruning and ...

\* Xg boost = ML algorithm

→ Extreme gradient boosting

Consider,

<u>Salary</u>	<u>Credit</u>	<u>Approval</u>	$y$	$\hat{y}$	$(y - \hat{y})$	$\frac{(y - \hat{y})}{B_1}$
$\leq 50K$	B	0	0	0.5	-0.5	-0.5
$\leq 50K$	B <sub>1</sub>	1	1	0.5	0.5	0.5
$\geq 50$	B <sub>1</sub>	1	1	0.5	0.5	0.5
$\geq 50K$	B	0	0	0.5	-0.5	-0.5
$\geq 50K$	B <sub>1</sub>	1	1	0.5	0.5	0.5
<del><math>\geq 50K</math></del>	<del>B<sub>1</sub></del>	<del>1</del>	<del>1</del>	<del>0.5</del>	<del>0.5</del>	<del>0.5</del>
$\leq 50K$	N	0	0	0.5	-0.5	-0.5

Steps

Step 1: Create a base model



← Base model is not biased to anything.

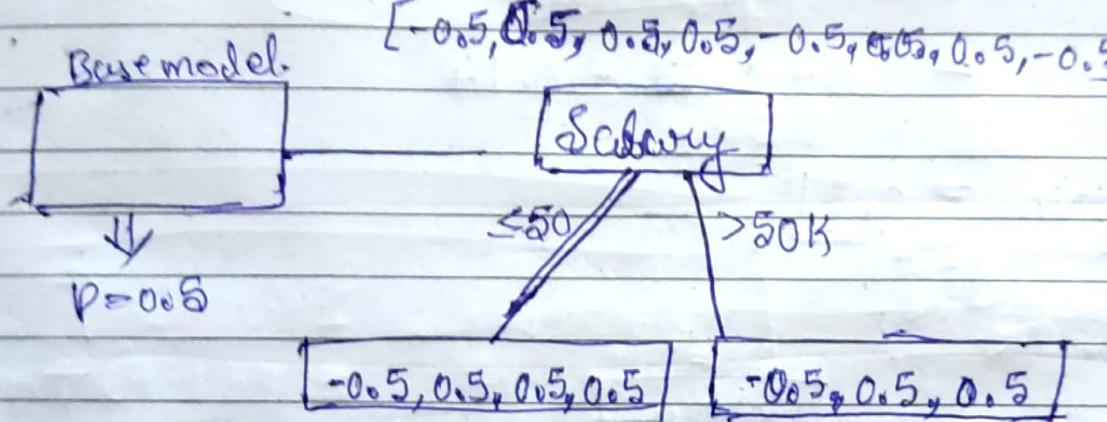
$$\downarrow \\ P = 0.5$$

$$\hat{y} = 0.5$$

Step 2: Construct decision tree with root

\* Salary and Credit are o/p and coefficient  $B_1$  is o/p  
 Calculate similarity weight  $\rightarrow$  calculate gain

~~step 3 and step 4 tells which feature to be selected as root~~



\* Similarity weight [LC] =  $\frac{(\sum \text{Residual})^2}{\sum P(1-P)}$

$$= \frac{(-0.5 + 0.5 + 0.5 - 0.5)^2}{[0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5)]}$$

Similarity weight [LC] = 0

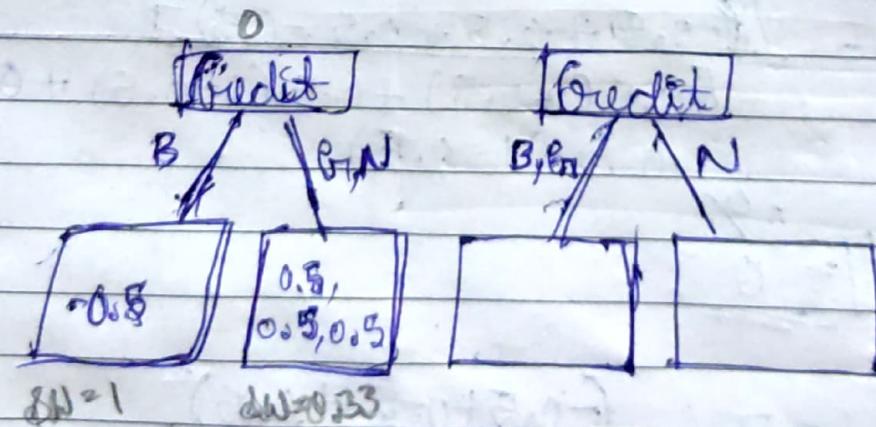
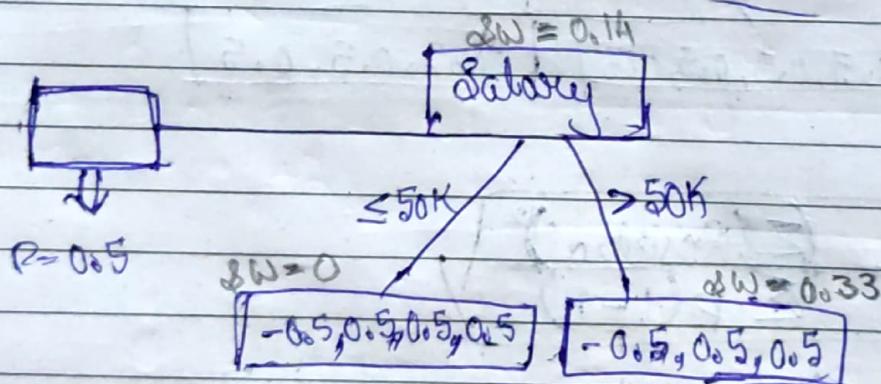
\* Similarity weight [RC] =  $\frac{(-0.5 + 0.5 + 0.5)^2}{[0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5)]}$

$$= \frac{0.25}{0.75} = \frac{1}{3}$$

Similarity weight [RC] = 0.33

$$\begin{aligned}
 & \text{symmetry weight [root]} = \frac{0.5 + 0.5 - 0.5 - 0.5}{0.5 + 0.5} = 0.75 \\
 & = \frac{1}{\sqrt{2}} = \underline{\underline{0.707}}
 \end{aligned}$$

$$\begin{aligned}
 & \text{gain} = \left[ \text{left.} + \text{right} - \text{root} \right] \\
 & \quad \quad \quad \text{child child node} \\
 & = 0 + 0.33 - 0.14 \\
 & \text{gain} = \underline{\underline{0.21}}
 \end{aligned}$$



$$* \text{ Similarity } [AC] = \frac{0.25}{0.25} = 1$$

$$\text{Similarity } [BC] = \frac{0.25}{0.75} = \frac{1}{3} = \underline{\underline{0.33}}$$

$$\text{Gain} \rightarrow 1 + 0.33 - 0 = \underline{\underline{1.33}}$$

\* Similarity, other child.

\* Now, Post any new data

*when passes through base model*

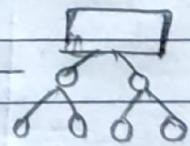
$$\log(\text{odds}) = \log\left(\frac{P}{1-P}\right)$$

$$\Rightarrow \log(\text{odds}) = \log\left(\frac{0.5}{1-0.5}\right) = \underline{\underline{0}}$$

~~Best fit~~

$$\rightarrow \boxed{\quad} \quad o/p = 0$$

$$\underline{\underline{0}}$$



For new record  
↳ 50K, B, 0

① New Test Data

$$= \Gamma(0 + \alpha(1))$$

*Base model o/p*

*Learning rate*

*Similarity weight from DT for record  
if  $\alpha = 0.1$*

*Logistic activation function*

$$= \Gamma(0.1)$$

$$= \frac{1}{1 - e^{-0.1}} = \underline{\underline{0.52}}$$

$$P(z) = \frac{1}{1 - e^{-z}}$$

$$\hat{y}$$

$\text{2nd record}$

②  $T \leq 50, \cdot B_1$

$$O/P = \sqrt{0 + d(0.33)}$$

$$= \sqrt{0 + 0.1 \times 0.33}$$

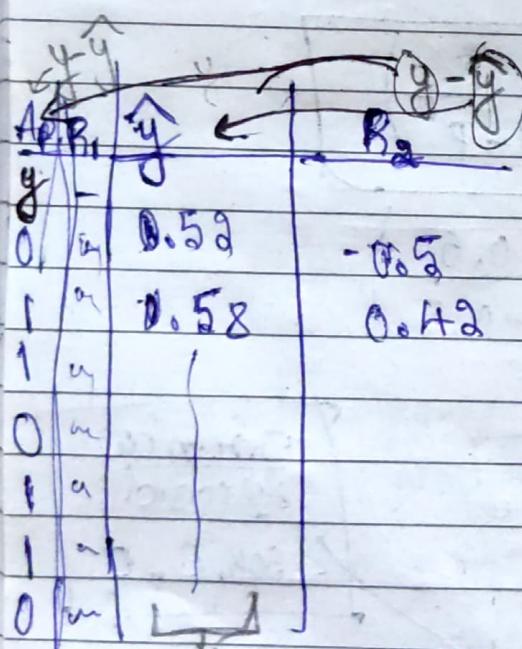
$$= \sqrt{0.033}$$

$$\Rightarrow \frac{1}{1 + e^{-0.033}}$$

$y=0.5$  (base model)

for 2nd record

$$O/P = \underline{\underline{0.508}}$$



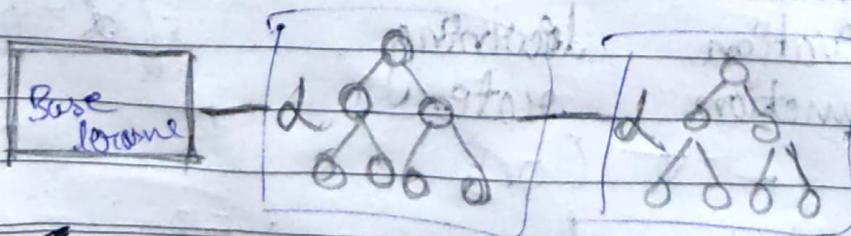
first model  
O/P

$$\text{similarity weight} = \frac{\sum (\text{Residual})^2}{\sum P(1-P) + \lambda}$$

hyperparam

What if when we have to set up split do we have to keep restriction to similarity weight.  $\Rightarrow$  "Cover Value"

$$\text{Cover Value} \Rightarrow \sum P_i(1-P_i)$$



$$O/P = \sqrt{[\text{Base learner} + d_1(\Delta T_1) + d_2(\Delta T_2) + \dots + d_n(\Delta T_n)]}$$