

# Bagging

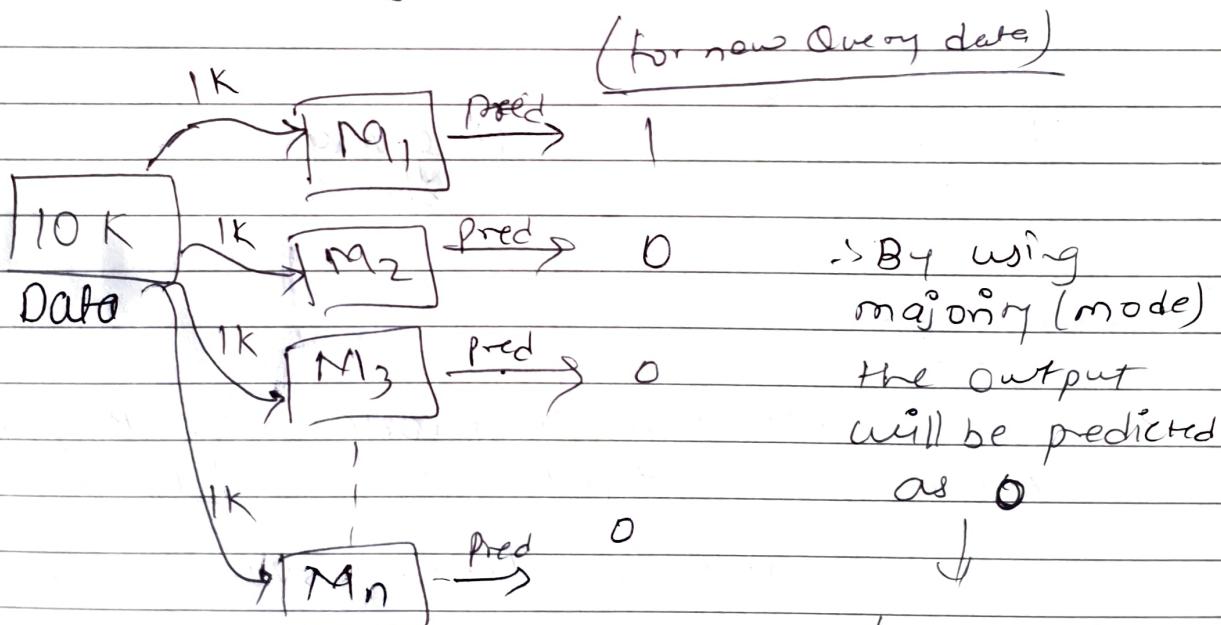
Bootstrapping

Aggregation

technique where population is being drawn from random sample.

All the base models ~~is~~ used for Bagging is ~~same~~ same.

for e.g. If I have 10 models then all 10 models should be having similar kind.



all these 1000 rows (data) given to each model ~~using~~ sampling will be diff. while training

Known as Bootstrapping

known as aggregation.

## why Bagging?

→ for better result,

[LB & LV]

there is negative correlation bet' bias & variance.

$$\text{bias} \propto \frac{1}{\text{variance}}$$

→ Ideally getting LB & LV is not feasible all the time due to their relation.

∴ Either we get LB HV  
or HB LV

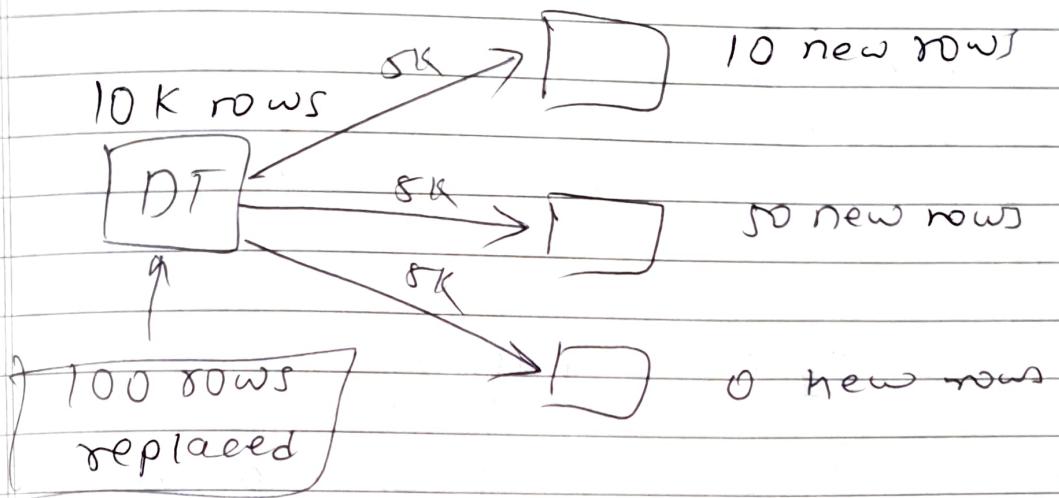
Hence, By using Bagging,  
LB & LV can be achieved.

## Explanation

while selecting a model for Bagging  
consider models with LB & HV

such as Decision Tree, SVM, KNN

Note: LB & HV will have overfitting issues



the above diag. says if 100 new rows are added/replaced then due to bootstrapping it is always the case that all newly added rows will not come under training data of each base model.

Hence biased & variance both will be low in such cases.

Or  
consistent results can be obtained.

## When to Use Bagging?

→ whenever you have a project or dealing with LB HV data then in such cases bagging will be useful.

### types of Bagging-

① Pasting: Row sampling same as bagging but without replacement.

② Random subspaces:- column sampling with or without replacement.

③

```
df = pd.read_csv('iris.csv')
df = df.sample(10)
df.sample(2, replace=True, axis=1)
```

↓                    ↓                    ↓  
 $(\text{NO. OF } \text{columns})$        $(\text{with replacement})$        $(\text{representative column sampling})$

④ Random patches :- Row sampling & column sampling

```
df.sample(8, replace=True).sample(2, replace
= True, axis=1)
```

# Bagging classifier

## Keypoints Related to bagging

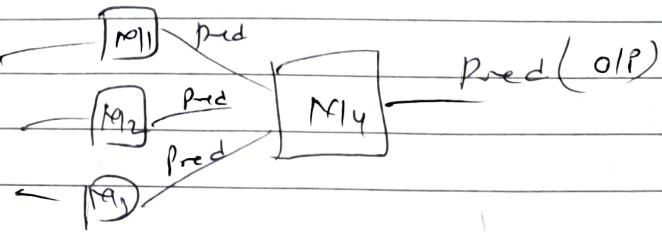
- ① Bagging generally gives better results than pasting
- ② Good results come around the 25% to 50% row sampling mark
- ③ Random patches and subsamples should be used while dealing with high dimensional data
- ④ To find the correct hyperparameter values we can do gridsearchcv / randomsearchcv

## Stacking

It extends ideas of max voting.

### Steps :-

- ① Train your base models.
- ② Predict the results using base model.
- ③ based on the prediction of base model, take that o/p and use it to train the meta model.
- ④



## Stacking vs Bagging (Boosting)

- |   |                                  |                               |
|---|----------------------------------|-------------------------------|
| ① | diff base models<br>can be used. | Same base models<br>are used. |
| ② | New training                     | No new training               |

### Drawback of stacking

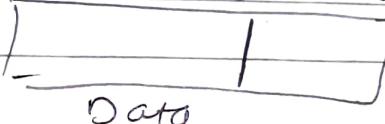
Overfitting.

→ Reason is that we are using same data for the ~~the~~ training of base model & then on the basis of same data only we are training the meta model too. Hence if there is overfitting while training the base model (which is always the case) then predicted result from meta model will also have overfitting.

### How to overcome?

- ① K - Fold method → known as stacking.

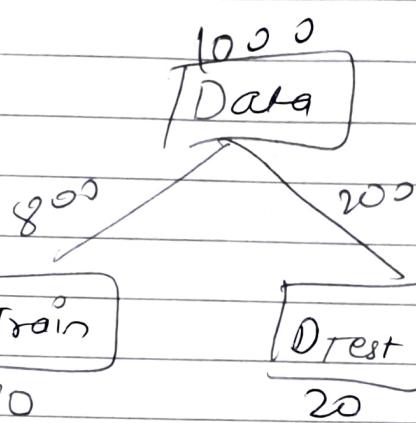
- ② Hold out method → known as blending



## blending (hold out)

Cgpa | iq | package

7	70	2.5
9	101	5.1
7.9	120	6.0
-	-	.



1000 Rows train

LR

DT

KNN

SVM / RFR

80

20

LR will be tested  
on Dvalidation &  
Dtrain will be stored as  
new dataset.  
→ similarly for all the  
base models, prediction  
will be done on Dvalidation

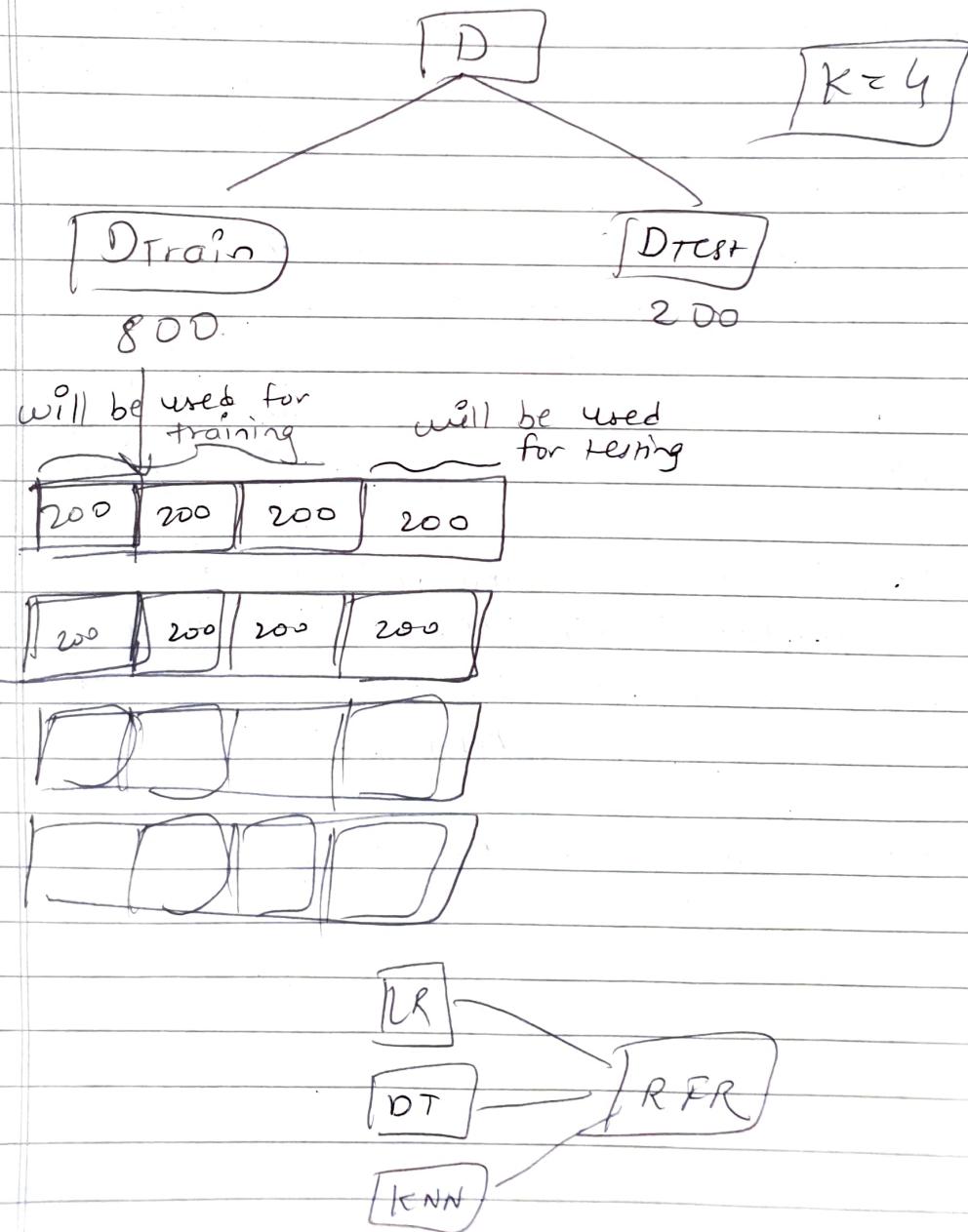
- STEP:
- ① train 3 base models (3 trained on Dtrain)
  - ② Do the prediction on m1, m2 + m3 on Dvalidation
  - ③ You form new dataset (using prediction of each base model)
  - ④ Train the meta model on Dtest (200 data)

## ② Stacking (K-fold)

→ K is how many equal parts of - training dataset.

Generally the value of K = 5 // K = 10.

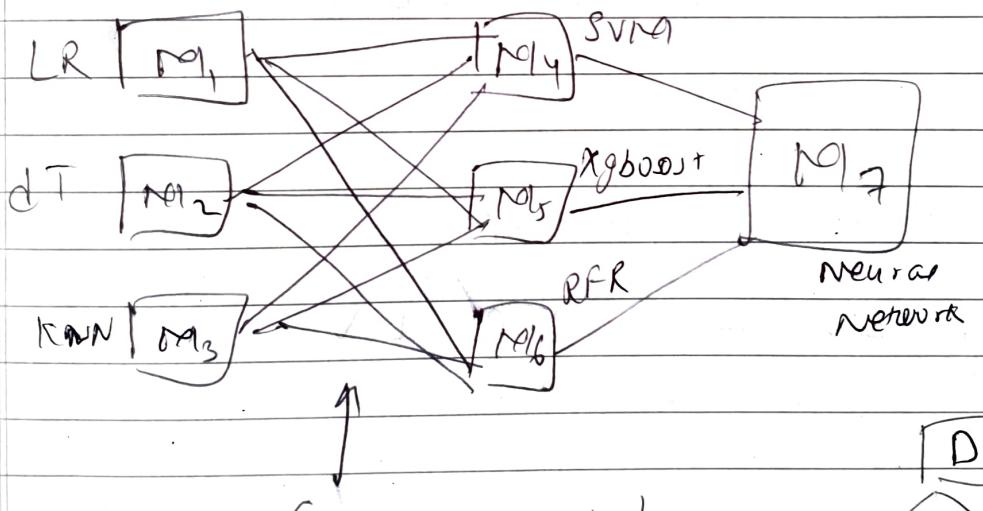
→ for K more than 10, computation will be more hence not suggested.



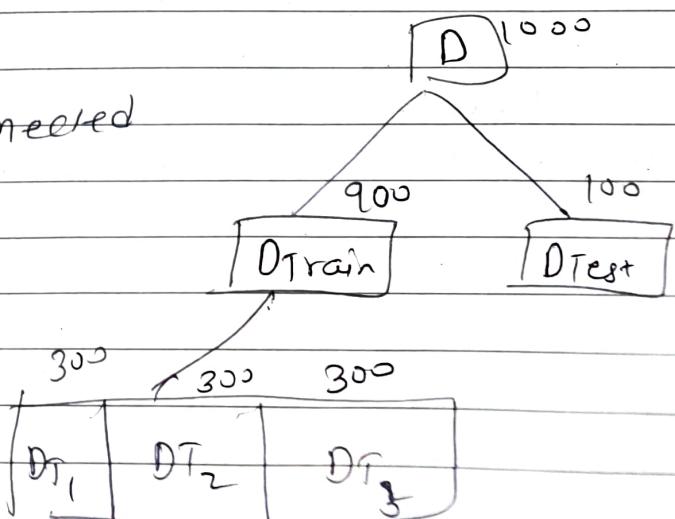
## Steps:

- ① while training for each base model  
 output will be  $(800, 4)$   
 and i.e. each base model will be trained K Times (4 times in this case)
- ② and then this predicted data will be given as training data (input data) to Meta model.
- ③ Retrain on  $X_{train}$  &  $y_{train}$  to train base model.

### Multi-layer stacking



fully connected



①  $M_1, M_2$  &  $M_3$  trained on data  
 $DT_1$

② prediction of  $M_1, M_2$  &  $M_3$  (i.e. the dataset which we get after prediction) will be given as training data to  $M_4, M_5$  &  $M_6$ .

③ prediction of  $M_4, M_5$  &  $M_6$  will be used to train  $M_7$

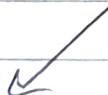
### Boosting

Bagging is nothing but conversion of LB & HV model into LB & LV

Similarly,

boosting is used to convert

HB LV  $\longrightarrow$  LB LV



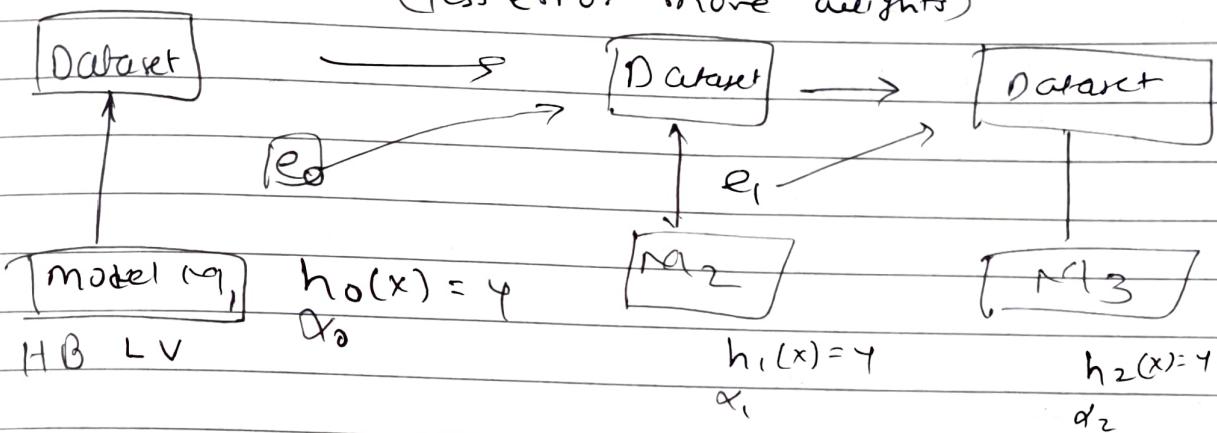
Decision tree  
of less depth.

$(x, y)$

$\alpha$  = weight assign to model  
based on error.

PAGE NO. / /  
DATE: / /

(less error more weights)



The points on which the error is coming will be boosted in the next step.

- No sampling is used in boosting
- Depending on the dataset, boosting stages are defined. ( $n$ -estimator = 10/100/1000/...)

$$\Rightarrow f = \alpha_0 h_0(x) + \alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots + \alpha_n h_n(x).$$

### Boosting technique

- ① Gradient boosting
- ② Xg boost (extreme boost)
- ③ Adaboost (adaptive boosting)

# Boosting techniques

## Ada boost

Multiple weak learners are added together to form strong powerful model.

### ① Weak learners:-

→ Generally their accuracies are less than 50% or just over 50%.

### ② Decision stumps:-

type of weak learner with maximum depth as 1.

for e.g.

$x_1$	$x_2$	$y$	Cgpa
6.4	79	1	+
4.9	81	0	-
9.1	102	1	+

→ Decision stumps gives better result in case of adaboost as weak learner.

### ③ $+1$ & $-1$

Binary classes of Adaboost.

$+1 \rightarrow +ve$

$-1 \rightarrow -ve$

There are no zero's (0) is used in classification.

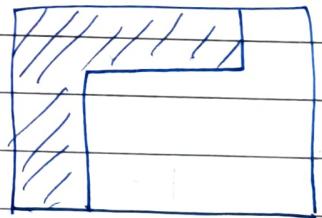
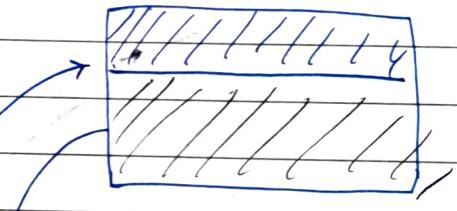
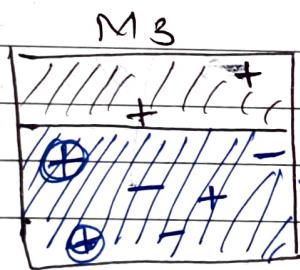
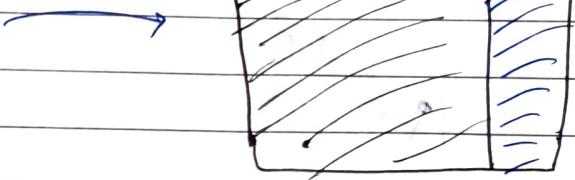
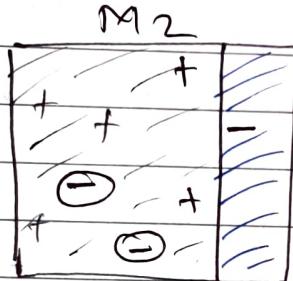
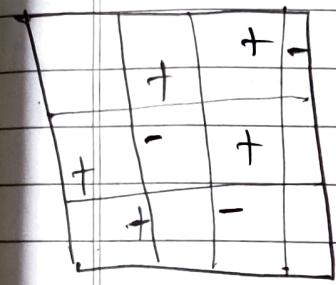
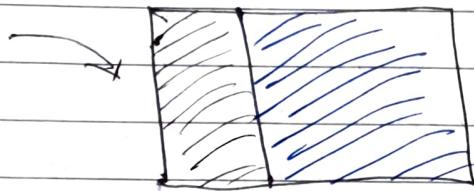
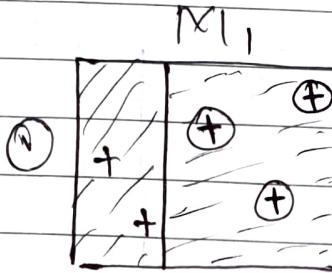
# Geometric Intuition

PAGE No. / /  
DATE: / /

\* Adaboost is a stagewise additive method

weak learners are added one by one

Adding weak learners all together



$$h(x) = \text{sign}(\alpha_1 h_1(x) + \alpha_2 h_2(x) + \alpha_3 h_3(x))$$

## Adaboost - Step by step

PAGE No.	1
DATE:	/ /

(Update)

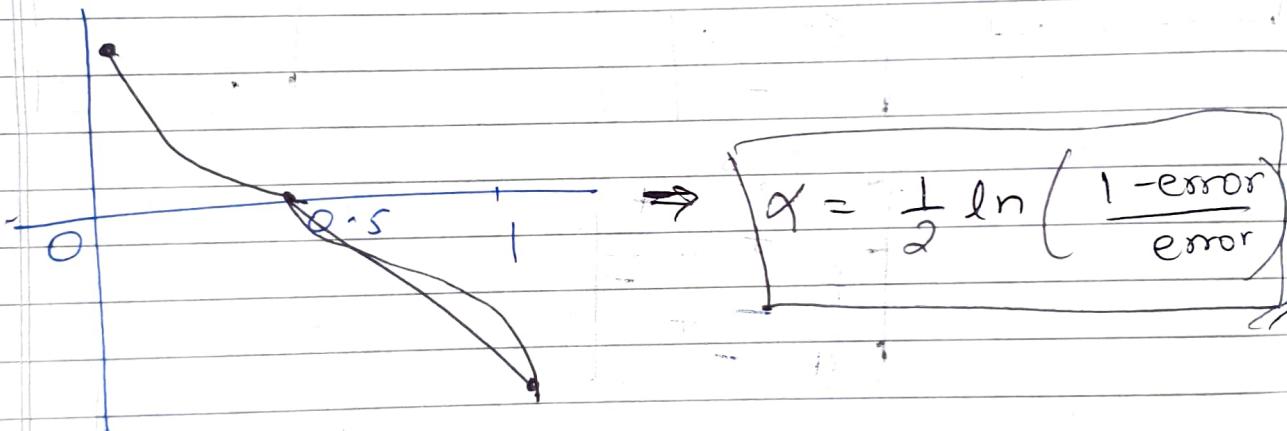
(Supersedes)

$x_1$	$x_2$	$y$	$y\text{-pred}$	weight	updated wt	normalized
3	7	1	1	0.2	0.16	0.166
2	9	0	1	0.2	0.24	0.25
1	4	1	0	0.2	0.24	0.25
9	8	0	0	0.2	0.16	0.166
3	7	0	0	0.2	0.16	0.166

## ~~Stage 1~~

Initial weight assign =  $\frac{1}{n}$  → no. of rows.

$$= \frac{1}{5} = 0.2$$



Now, how to calculate error of a model.

$\therefore$  error = sum of weight of all misclassified points.

: For M<sub>1</sub>

$$\text{error} = 0.2 + 0.2 \\ = 0.4$$

$$\therefore \alpha_1 = \frac{1}{2} \ln \left( \frac{1-0.4}{0.4} \right)$$

$$= \frac{1}{2} \ln \left( \frac{0.6}{0.2} \right)$$

$$\boxed{\alpha_1 = 0.20}$$

stage 2:-

Now,  $M_1$  has to inform  $M_2$  that some mistake has been done so be careful.

$\therefore$  the weights of misclassified points will be [increased]  $\rightarrow$  (boost)

and

the weights of correctly classified points will be reduced.

$\therefore$  For misclassified

$$\text{new-wt} = \text{curr-wt} \times e^{\alpha_1}$$

For correctly classified

$$\text{new wt.} = \text{curr wt.} \times e^{-\alpha_1}$$

$\therefore$  for previous table

$$\text{new-wt} = 0.2 \times e^{+0.2} = 0.24 \quad (\text{for misclassified})$$

$$\text{new-wt} = 0.2 \times e^{-0.2} = 0.16 \quad (\text{for correctly classified})$$

$x_1$	$x_2$	$y$	new-wt	range	
3	7	1	0.166	0 - 0.166	
2	9	0	0.25	0.166 - 0.416	
1	4	1	0.166 0.25	0.416 - 0.666	0.13
9	8	0	0.166	0.666 - 0.832	0.43
3	7	0	0.166	0.832 - 1.0	0.62

for upsampling

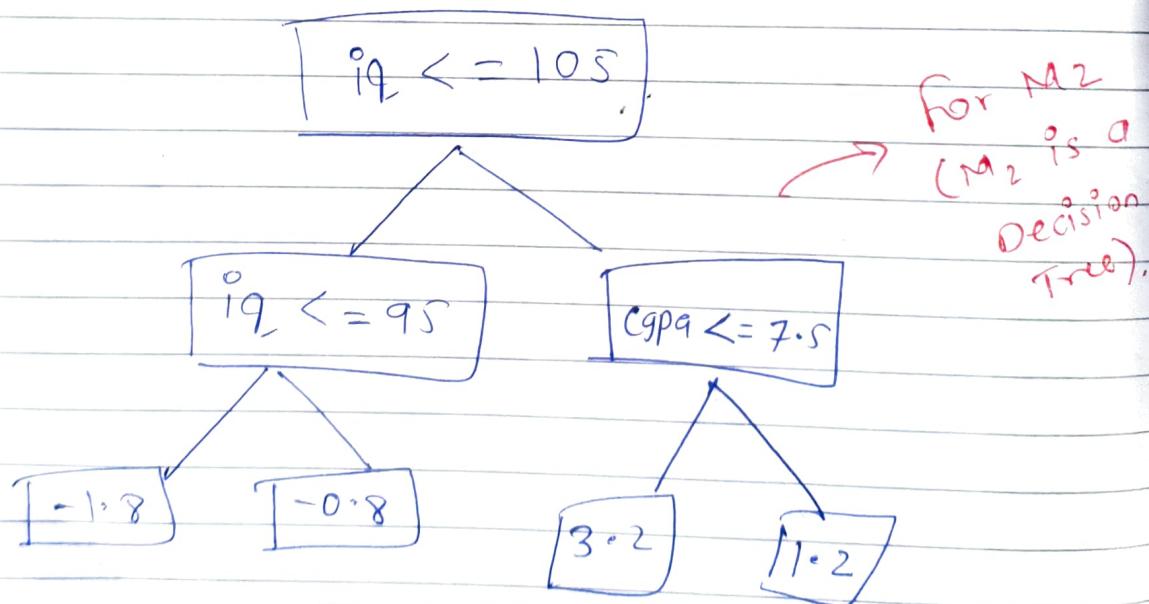
model will generate 5 random nos. (coz n=5)  
based on weightage and find one  
~~where~~ the rows for that each no.  
depending on range.

After that new dataset will be prepared  
where row no. 1, row no. 3 (3 times)  
and row no. 4 will be there.

# Gradient Boosting (Regression).

$x_1$	$x_2$	$y$	$M_1$	pseudo-residual	Pred-2	res-2	Pred 3	res-3
			Pred-1	$= \text{Actual} - \text{pred}$ $= y - \text{Pred-1}$				
90	8	3	4.8	-1.8	-1.8	-1.62		
100	7	4	4.8	-0.8	-0.8	-0.72		
110	6	8	4.8	3.2	3.2	2.88		
120	9	6	4.8	1.2	1.2	1.08		
80	5	3	4.8	-1.8	-1.8	-1.62		
			24					

Note:- In Gradient boosting, for 1<sup>st</sup> model in case of regression is calculating simply the mean of output.



$$y_{\text{pred}} = m_1 + LR \times m_2$$

$$= m_1 + 0.1 \times m_2$$

$$Res_2 = \text{Actual} - (\frac{\text{pred}_1}{3} + 0.1 \times \text{pred}_2)$$

for  $x_1, 3$

$$P_1 \leq 105$$

$$P_2 \leq 95$$

$$Capa \leq 75$$

$$-1.62 - 0.72$$

$$2.88$$

$$1.08$$

at this point,

prediction will be given as,

$$4\text{-pred} = m_1 + 0.1 \times M_2 + 0.1 \times m_3$$

# Gradient boosting

## (classification)

→ It is similar as logistic Regression.

(E)

	A engg. Likes	B Age	C fav. sub	D Likes offline lec?	Residual (actual - Pred) Using eqn ②	New Pred - Probability	New Residual
①	Yes	12	ML	Yes	0.3	0.9	0.1
②	Yes	87	AI	Yes	0.3	0.5	0.5
③	No	44	ML	No	-0.7	0.5	-0.5
④	Yes	19	DL	No	-0.7	0.1	-0.1
⑤	No	32	AI	Yes	0.3	0.9	0.1
⑥	No	14	ML	Yes	0.3	0.9	0.1

① Initial prediction for every individual is the log( odds).

$$\begin{aligned}
 &= \log \left( \frac{\text{total yes}}{\text{total no}} \right)_0 \\
 &= \log \left( \frac{4}{2} \right) \\
 &= 0.7 \text{ (approx)}
 \end{aligned}$$

just like logistic regression, to use the log( odds) for classification is to convert it to a probability.

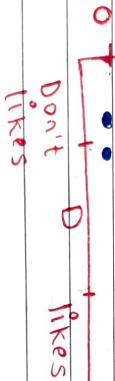
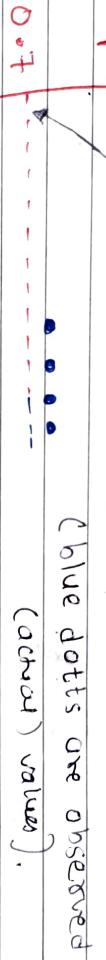
$$\begin{aligned}
 \therefore \text{probability of like offline lec} &= \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}} \\
 &= 0.7 \text{ (approx)}
 \end{aligned}$$

\* Since probability of log(odd) is greater than 0.5, we can classify everyone (every row) in the training dataset as someone who loves (likes) offline lectures.

Based on above conclusion (prediction) it is observed that predicted prediction is not correct.

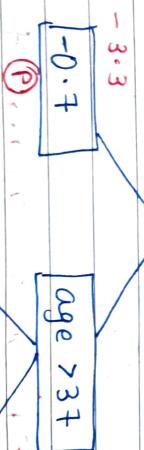
② Calculating residual = (Observed - Predicted)

Predicted values



③ Now, build a tree using column A, B, C to predict residuals.

$$\text{Fov\_sub} = \text{OL}$$



-1

$$[0.3, -0.7]$$

1.4

$$[0.3, 0.3, 0.3]$$

⑧

⑨

Note:- While creating decision tree, maximum depth can vary from 8 to 32.

Now,

for Gradient boosting for classification, most common transformation formula is,

$$\text{OLP Value} = \frac{\sum \text{residual}^2}{\sum [\text{previous probability} \times (1 - \text{prev-probability})]} - \textcircled{1}$$

∴ for Row 4 (for P) [from Decision tree]

$$\begin{aligned}\text{OLP value} &= \frac{-0.7}{0.7 \times (1-0.7)} = \textcircled{2} \\ &= -3.03\end{aligned}$$

for Q (from Decision tree)

$$\begin{aligned}&= \frac{0.3 + (-0.7)}{[0.7 \times (1-0.7)] [0.7 \times (1-0.7)]} \\ &= -1\end{aligned}$$

for R (from Decision tree)

$$\begin{aligned}&= \frac{0.3 + 0.3 + 0.3}{3 [0.7 \times (1-0.7)]} \\ &= 1.64\end{aligned}$$

(4) log(Odds) prediction = log(Odds) + \lambda \times (\text{OLP value from tree formed } \textcircled{1}) - \textcircled{2}

For e.g.:  
For Row 1 -  
 $\log(\text{Odds}) \text{ prediction} = 0.7 + 0.8 \times 1.4$

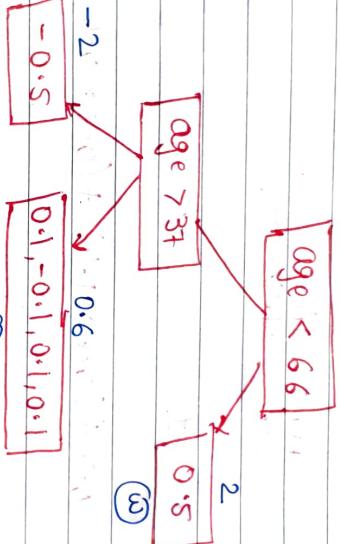
$\lambda = 0.8$  for illustration purpose. Generally  $\lambda$  is 0.1 ] ,

Now,

Converting new log(odd) prediction into probability

$$\text{probability } \log(\text{odds}) = \frac{e^{1.8}}{1+e^{1.8}} = 0.9$$

Now! based on column A, B, C and E (new-residual) we create a new tree.



④

⑤

Next,

Calculating o/p value for each leaf node ( $W, X, Y$ )

$$\textcircled{1} \quad \text{for } W: \quad = \frac{0.5}{0.5(1-0.5)} = \boxed{2}$$

$$\textcircled{2} \quad \text{for } X: \quad = \frac{-0.5}{0.5 \times (1-0.5)} = \boxed{-2}$$

$$\textcircled{3} \quad \text{for } Y: \quad = \frac{-0.1 + 0.1 + 0.1 + 0.1}{[0.9(1-0.9)][0.1 \times (1-0.1)][2 \times 0.9 \times (1-0.9)]} \\ = \boxed{0.6}$$

Now again,

$$\text{log(odds) prediction} = \text{log(odds)} + \lambda \times (\text{from tree 2})$$

\* Prediction \*

→ let's for an input,

likes	age	fav. sub.	likes offline recs?
engg.	25	AI	???

Yes      25      AI      ???

$$\text{log(odds)} \rightarrow \text{prediction} = 0.7 + (0.8 \times 1.4) + (0.8 \times 0.6) = 2.3$$

$$\downarrow \alpha_1 \times \text{OIP}$$

Value of  $\alpha_2 \times \text{OIP}$  value from

(R) From team 1  
(Y) From team 2

$$\therefore \text{log(odds)} = 2.3$$

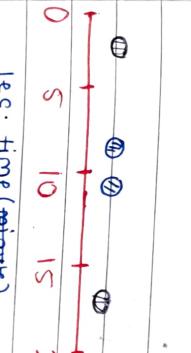
Now, use convert this log(odds) into a probability.

$$\text{probability} = \frac{e^{2.3}}{1 + e^{2.3}} = 0.9 \text{ (approx)}$$

Since, threshold value is 0.5  
∴ Class of the given I/P data will be Yes. ( $0.9 > 0.5$ )

# Xgboost [Extreme Boosting] $\Rightarrow$ Classification

→ Xgboost is mostly suitable for large, complicated datasets.



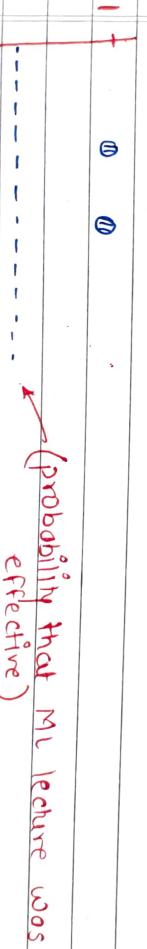
→ blue dots denotes that ML lecture was effective.

black dots represents that ML lecture was not effective.

Step 1: Making initial prediction. [This value can be anything but by default it is 0.5 for classification & Regression both]

Now,

Illustrating the initial prediction by adding a y-axis to our graph to represent the probability that ML lecture is effective.



Now,

Calculate residual = actual - predicted just to check whether how good the initial prediction is.

Now fit an xgboost tree to the residuals.

$$\text{Cleft.time} \\ \text{dosage} < 15$$



Next step is to calculate the similarity scores.

$$\text{similarity scores} = \frac{\sum (\text{Residual})_i^2}{2}$$

$$= \frac{\sum [\text{previous probability}_i \times (1 - (\text{prev-probability}_i))]^2}{\sum \text{residuals}^2}$$

for creating a tree

→ Each tree starts out as a single leaf.  
and all of the residuals go to the leaf.



→ Next, calculate the similarity score for the leaf.

$$= \frac{(-0.5 + 0.5 + 0.5 - 0.5)^2}{\text{Denominator}}$$

$$= 0$$

$$[-0.5, 0.5, 0.5, -0.5]$$

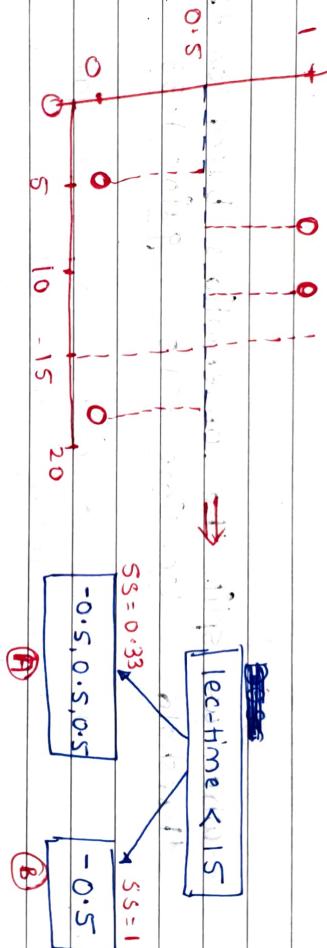
similarity = 0

Now, we need to decide if we can do a better job - clustering similar residuals if we split them into two groups.

for e.g.: let's start with this threshold,

lec-time < 15:

the reason we chose 15 because it is the avg. value b/w last 2 observations.



NOW, calculate the similarity score of the leaf on the left, and then on the right respectively.

\* Similarity of left leaf (A): :-

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

$$= \frac{[0.5 \times (1 - 0.5)][0.5 \times (1 - 0.5)]}{[0.5 \times (1 - 0.5)][0.5 \times (1 - 0.5)]} + \lambda$$

↳ [Note: Initial prediction was 0.5]

$$= 0.33$$

\* similarly, similarity score for right leaf (B): -

$$= \frac{(0.5)^2}{0.5 \times (1 - 0.5)} = 1 \text{ (when } \lambda = 0\text{)}$$

Now, for simplicity we'll put  $\lambda = 0$  but it is known that  $\lambda$  reduces the similarity score which makes leaves easier to find.

After calculating the similarity score, we need to calculate the Gain Information gain.

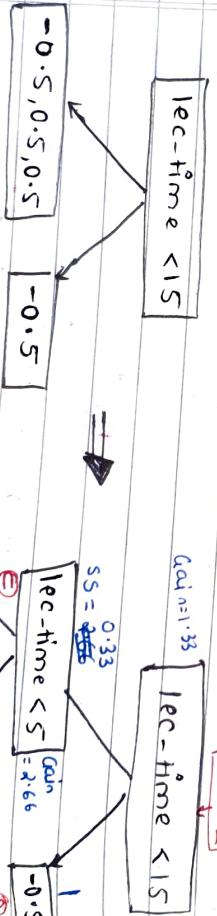
$$\text{Gain} = \text{Left Similarity} + \text{Right Similarity} - \text{Root Similarity}$$

$$\therefore \text{Gain of tree formed} = 0.33 + 1 - 0 \\ = 0.33$$

∴ When we split the observations based on the threshold lec-time < 15, gain = 1.33

\* \* \* Since, No other threshold will give gain greater than lec-time < 15, Hence lec-time < 15 will be the first branch of the tree.

Tree 2



This is because

$$(SS = 3.66) \quad (A)$$

(B)

gain than the lec-time < 10.  
 $(SS = 0.66)$

Note: Just for simplicity we are not splitting any leaf further.

i.e. limited the tree as max level = 2

\* The minimum number of residuals in each leaf is determined by cover.

Cover: Cover is defined as diff. bet'n denominator of similarity score and N.

$$\therefore \text{cover} = \sum [(\text{prev.-probability}_i) \times (1 - \text{prev.-probability}_i)] - N$$

By default, the minimum value for cover is +1.

Now, for Xgboost classification, cover depends on the previously predicted probability of each residual in a leaf.

for e.g. in tree 2, for leaf node A,

$$\text{cover} = 0.5 \times (1 - 0.5) = 0.25$$

Similarly, in tree 2, cover for leaf node B is,  $\min$

$$\text{cover} = 0.5 \times (1 - 0.5) + 0.5 \times (1 - 0.5) = 0.5$$

since, default value for cover is 1, xgboost would not allow the leaf A & B.

Hence it will be pruned.

Ultimately, if we use the default minimum value for cover as 1, then we would be left with the root, and xgboost requires trees to be larger than just the root.

Hence, in order to have a tree, we'll assume the minimum value for cover to be zero(0).

In program:-  $\boxed{\text{cover}(\text{min\_child\_weight}) = 0}$

### Condition for pruning:-

$\text{cover} - \gamma = \begin{cases} \text{if positive, then do not prune} \\ \text{if negative, then prune.} \end{cases}$

( $\gamma$  = tree complexity parameter)

more the no. of  $\gamma$ , less will be the different (negative result). Hence causes to prune the branches.

Hence value of  $\gamma$  and  $\lambda$  will be selected as zero in this example for simplification.

For now,

regardless of  $\lambda$  and  $\gamma$ , let's assume that this is the tree (tree 2) ~~is the~~ we are working with.

\* Determining OIP values for the leaf leaves. (Tree(2))

$$\text{OIP values} = \frac{\sum \text{residuals}}{\lambda}$$

$$\sum [\text{prev-probability}_i \times (1 - \text{prev-probability}_i)] + \lambda$$

$\lambda$  = Regularization parameter

$$\begin{aligned} &\therefore \text{OIP value for leaf A (in tree(2))} \\ &= -0.5 \\ &= 0.5 \times (1 - 0.5) + 0 = -2 \end{aligned}$$

(considering  $\lambda=0$  in this example)

$\therefore \text{lect-time} < 15$

$$\begin{array}{c} \text{lect-time} < 5 \\ \diagdown \quad \diagup \\ [-0.5] \qquad [-0.5] \end{array}$$

$$\text{OIP} = -2$$

$$\begin{array}{c} [-0.5] \qquad [0.5, 0.5] \\ \diagup \quad \diagdown \\ \text{OIP} = -2 \end{array}$$

$$\text{OIP} = -2$$

∴ 1st tree is completed after calculating the OIP value of each leaf.

Now, based on the tree formed, new prediction will be done.

→ Just like other boosting methods, Xgboost classification makes new predictions by starting with the initial prediction.

Now, calculating log(odd).

$$\begin{aligned}\therefore \text{log(odd)} &= \log\left(\frac{P}{1-P}\right) \\ &= \log\left(\frac{0.5}{1-0.5}\right) - [\because \text{initial prediction was } 0.5] \\ &= \log(1) \\ &= 0\end{aligned}$$

→ For Xgboost classification, learning rate ( $\epsilon$ ) is used.

The default value for  $\epsilon$  is 0.3.

$$\therefore \text{prediction} = \text{log(odd)} + \epsilon \cdot x \quad \begin{matrix} (\text{OIP values}) \\ (\text{from the tree formed}) \end{matrix}$$

Let's calculate for lect-time = 2

$$\begin{aligned}&= \text{log(odd)} \text{ prediction} \\ &= 0 + 0.3 \times (-2) \\ &= -0.6\end{aligned}$$

→ To convert log(odd) value into probability, we use logistic function (sigmoid function).

$$e^{\text{log(odd)}}$$

$$\therefore \text{log(odd) probability} = \frac{e^{\text{log(odd)}}}{1 + e^{\text{log(odd)}}}$$

$$= \frac{e^{-0.6}}{1 + e^{-0.6}} \\ = \boxed{0.35}$$

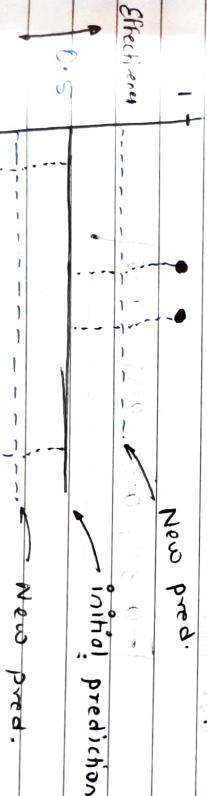
similarly, we'll calculate the log(odd) probability for each observation.

Calculating log(odd) prediction for lect-time = 8

$$= 0 + 0.3 \times 2$$

$$= 0.6$$

$$\therefore \text{log(odd) probability} = \frac{e^{0.6}}{1 + e^{0.6}} = \boxed{0.65}$$



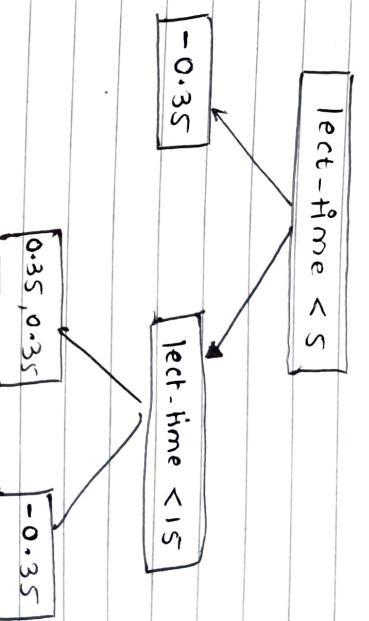
(let-time)  $\rightarrow$

Conclusion: By looking at above graph, it is observed that newly calculated residual is smaller for each observation than the initially calculated value of residual.

Since, less the residual more is the accuracy.

$\therefore$  Xgboost is enhancing step by step.

NOW, based on the new residual we can fit a new tree.



Note:-

While building the new tree, calculating the similarity score is more interesting because previous probabilities are no longer same for

$\therefore$  root will look like,

$$[-0.35, 0.35, 0.35, -0.35]$$

OP value for root,

$$\sum \text{residual}^2$$

$$= [0.35 \times (1-0.35)] + [0.65 \times (1-0.65)] + [0.65 \times (1-0.65)] + [0.35 \times (1-0.35)]$$

# Random Forest

PAGE N  
DATE:

Random forest with gradient boosting or Xgboosting is one of the most widely used algorithm in ML.

Applicable for classification & Regression.

→ Prerequisites

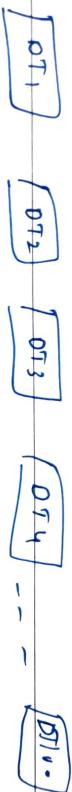
- ① Decision tree
  - ② Bagging Techniques

Q. Why Random Forest?  
↳ multiple ~~trees~~ Decision trees are trained together

It is a bagging technique only. But when decision tree is used as a base model in bagging then it becomes Random forest.

$\therefore$  Random forest is specialization of bagging technique.

Dataset = 1000 Rows



$\rightarrow$  Row sampling }  
 (feature sampling)  $\rightarrow$  column sampling } with / without  
 $\rightarrow$  combination }  
 $\quad \quad \quad$  Replacements }

## Bagging vs Random Forest

### Bagging

RF

- ① Diff. algorithms can be used. ex. KNN, DT, SVM as a base model.

- ② Tree level column sampling / Node level sampling

on RF is ~~be~~ having better performance than Bagging.