

# Decision Tree

①

└ ID3

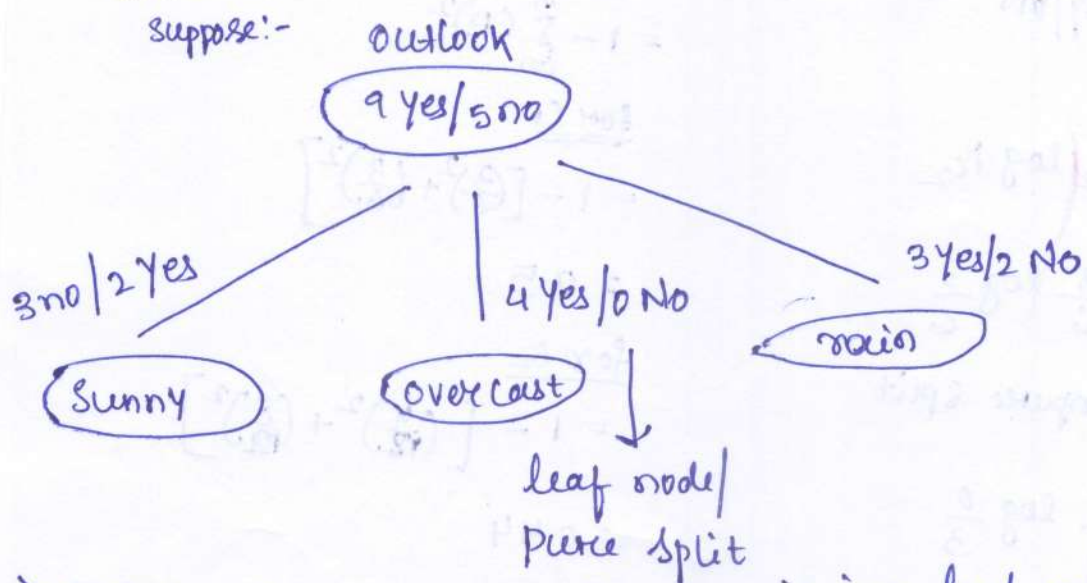
└ CART (Classification & Regression Tree)

What is a leaf node?

⇒ No more splitting happens if leaf node occurs

Ex:-

Suppose:-



→ Splitting happens until/unless it is a leaf node.

How application will know whether it is a leaf node or not

→ Purity of a leaf node is checked by using

→ Entropy

→ Gini Index

## Entropy

→ used for small data set.

→ It ranges between (0-1)

Binary classification

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

Multi class classification

$$H(S) = -P_{C1} \log P_{C1} - P_{C2} \log P_{C2}$$

## Gini Index

→ for large dataset

→ It ranges between (0-0.5)

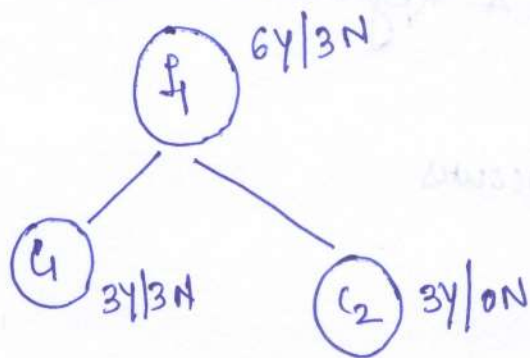
$$G.I = 1 - \sum_{i=1}^n (P_i)^2$$

for Binary

$$= 1 - \sum_{i=1}^n ((P_+)^2 + (P_-)^2)$$

## Entropy

Ex:-



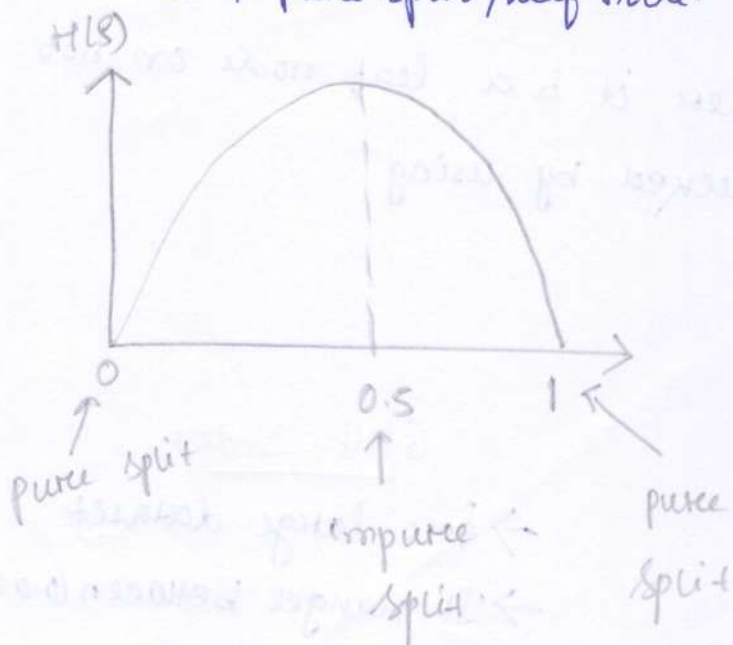
$$H(S) = -P_{C1} \log P_{C1} - P_{C2} \log P_{C2}$$

$$H(C1) = -\frac{3}{6} \log \frac{3}{6} - \frac{3}{6} \log \frac{3}{6}$$

$$= 1 \rightarrow \text{very Impure split}$$

$$H(C2) = -\frac{3}{0} \log \frac{3}{0} - \frac{0}{3} \log \frac{0}{3}$$

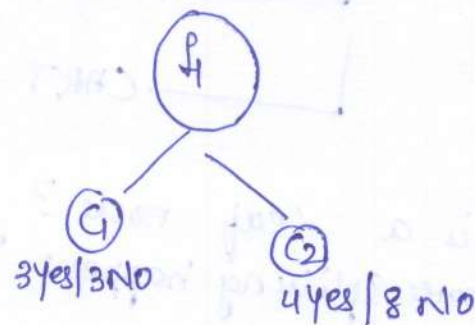
$$= 0 \rightarrow \text{pure split / leaf node.}$$



## G.I

②

Ex:-



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

for C1

$$= 1 - \left[ \left(\frac{3}{6}\right)^2 + \left(\frac{3}{6}\right)^2 \right]$$

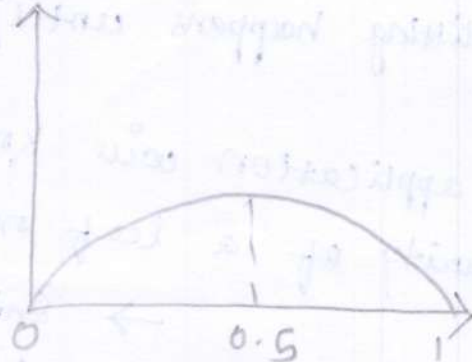
$$= 0.5$$

for C2

$$= 1 - \left[ \left(\frac{4}{12}\right)^2 + \left(\frac{8}{12}\right)^2 \right]$$

$$= 0.44$$

H(S)





How Features

which features to take the split?

→ It is decided by Information Gain

Information Gain

$$\text{Gain}(S, f_1) = H(S) - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} H(S_v)$$

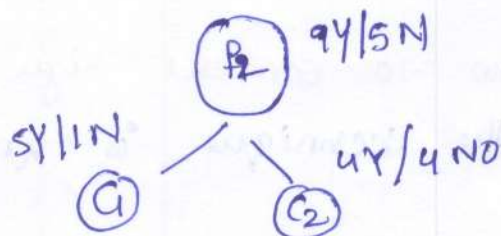
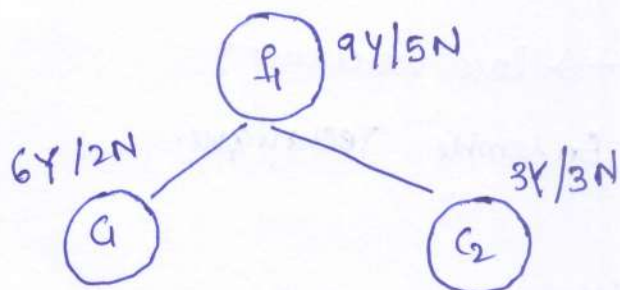
$H(S)$  = entropy / G.I

$H(S_v)$  = entropy of child node

$|S|$  = Total nodes

$|S_v|$  = child nodes.

Ex:-



entropy of root

$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$\approx 0.94$$

$$H(C_1) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}$$

$$= 0.81$$

$$H(C_2) = 1$$

$$\text{Gain}(S, f_1) = 0.94 - \left[ \left( \frac{8}{14} \right) \times 0.81 + \left( \frac{6}{14} \right) \times 1 \right]$$

$$= 0.049$$

$$H(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.937825$$

$$H(C_1) = -\frac{5}{6} \log_2 \frac{5}{6} - \frac{1}{6} \log_2 \frac{1}{6} = 0.6497$$

$$H(C_2) = 1$$

$$\text{Gain}(S, f_2) = 0.937825 - \left[ 0.6497 \times \frac{6}{14} + 1 \times \frac{8}{14} \right]$$

$$= 0.088425$$

$$I.G(f_2) > I.G(f_1)$$

∴

So we use feature-2 for splitting

Why over-fitting happens in Decision Tree?

(4)

→ If decision tree is decision to its full depth

↳ Training accuracy  $\uparrow\uparrow$  (low bias)

↳ Testing "  $\downarrow\downarrow$  (High Variance)

How to fix overfitting?

Pre-pruning

→ Initially max depth is chosen

Post-pruning

→ First the DT is built to depth then cutting is done.

How to convert high variance  $\rightarrow$  low variance?

→ The technique is known as Ensemble Technique.

Ensemble Technique

Boosting

→ Each algorithm are parallel

→ Majority voting classifier

→ Avg o/p in Regressor

Ex:-

Random Forest classifier  
- Regressor.

Bagging

→ Sequential Algorithm is applied.

Ex:-

Adaboost, Catboost, Gradient boosting, XG-Boost.



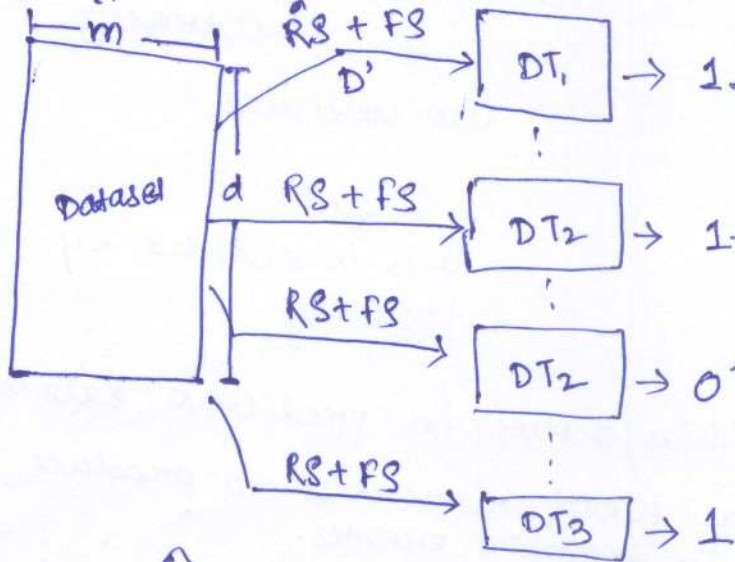
## Random Forest classifier

$$D' < D$$

→ no of columns

RS + FS

D'



RS → Row Sampling

FS → Feature Sampling

⑤

→ Row sampling with Replacement occurs in the step-① i.e. new Rows of dataset is picked there might be few repetition.

To aggregate the majority Vote is considered.

$$o/p = 1$$

↑  
When Test data is given

↑  
prediction Takes place.

## Advantages

① To make high variance → low variance we are using multiple Decision Tree in parallel and doing majority voting

② If we are changing data in the dataset it will not affect much the accuracy, bcz the data are splitted properly within the decision tree not causing major impact.

## If it is a Regression

→ we take mean or median of the output

How many Decision Tree to be used?

• It can be decided by hyper-parameter Tuning.

## Bagging

① High variance + low bias



Variance is reduced by Randomization + Aggregation

## Boosting

②

High bias (low training accuracy) +

low variance

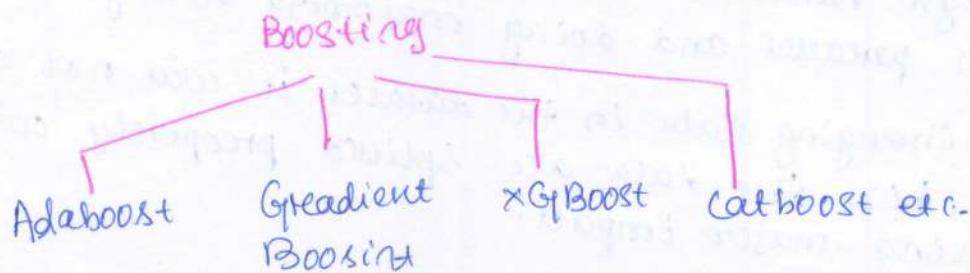


bias is reduced by boosting

Boosting is a method to reduce bias/error in predictive data analysis

→ It combines set of sequential weak learners to produce a Strong learner to minimize the training errors.

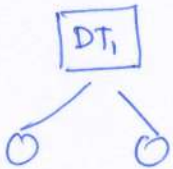
→ In boosting, a random sample of data is selected, fitted with a model and then trained sequentially i.e. each model tries to compensate for weakness of its predecessor with each iteration.



## Adaboost

→ Combines weak learners <sup>in sequential</sup> to make a Strong learner.

Stump = DT with one level & 2 binary o/p.



→ Each Stump is made by taking the previous Stumps mistakes into account.



<u><math>f_1</math></u>	<u><math>f_2</math></u>	<u><math>f_3</math></u>	<u>B/P</u>	<u>old weights</u>	<u>updated weights</u>	<u>New sample weight</u>
			Yes	$\frac{1}{7}$	0.058	$\frac{0.058}{0.697} = 0.0832$
			No	$\frac{1}{7}$	0.058	0.0832
			-	$\frac{1}{7}$	0.058	0.0832
			-	$\frac{1}{7}$	0.349	0.5007
			-	$\frac{1}{7}$	0.058	0.0832
			-	$\frac{1}{7}$	0.058	0.0832
			-	$\frac{1}{7}$	0.058	0.0832
					0.697	0.9921

Step 1  $\Rightarrow$  Equal weights are assigned i.e.  $= \frac{1}{N} = \frac{1}{7}$

Step 2  $\Rightarrow$  Then we will create decision stump for each feature and calculate G.I of each tree. Let tree with lowest G.I will be our first stump.

Let 4<sup>th</sup> one is our 1st Stump.

Step 3 we will calculate the "Amount of say" or "Importance" or "influence" i.e. performance of stump.

$$\frac{1}{2} \log_e \left( \frac{1-TE}{TE} \right) = \frac{1}{2} \log_e \left( \frac{1-1/7}{1/7} \right) \quad \therefore \text{Let one record is at fault so error} = \frac{1}{7}$$

$$\approx 0.895$$

Step 4 :- Weights need to be updated i.e. decrease the weight of correct record, increase the weight of incorrect record.

$$\begin{aligned} \text{For correct record} &= \text{Weight} \times e^{-I_S} \\ &= \frac{1}{7} \times e^{-0.895} = 0.058 \end{aligned}$$

$$\begin{aligned} \text{For wrong record} &= \text{Weight} \times e^{I_S} \\ &= \frac{1}{7} \times e^{0.895} = 0.349 \end{aligned}$$

We know that Total sum of updated weights must be equal to 1 but we are getting 0.697. To bring this sum equal to 1 we need to normalize these weights by dividing all the weights by the total sum of updated weights.



Then we will do bucketing, i.e. max bucket size is for wrong records i.e. the wrong predicted is passed to next decision tree. So that training of weak learners will be done on the next step.

## Gradient Boosting

→ The main idea behind this algorithm is to build models sequentially and these subsequent model try to reduce the error of the previous model.

Ex:-

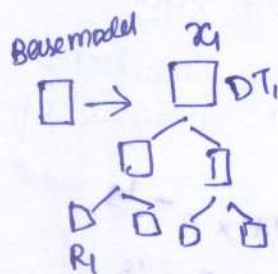
Exp	Degree	Salary	(predicted value) $\hat{y}$	( $\hat{y}_1$ ) $y - \hat{y}$	let suppose $R_2$	
2	BE	50K	75	-25	-23	$75 + (-23)$
3	masters	70K	75	-5	-3	$= 52$ (overfitting)
5	masters	80K	75	5	3	$75 + \alpha(R_2)$
6	PHD	100K	75	25	20	$75 + 0.1(-23)$ $= 73.7$

Step 1:- Compute the Base model to get single o/p.

$$\text{avg} = \frac{50 + 70 + 80 + 100}{4} = 75$$

Step 2:- Compute Residuals Error / pseudo Residuals ( $R_1$ )  
actual - prediction

Step 3:- construct DT  
 $\{x_i, R_i\}$



$$h_1(x) = \text{o/p of } DT_1$$

Step 4:- let we pass independent features to DT then  $R_2$  error comes

Here we say, predicted -  $R_2 = 52$  it is close to  $y$  but overfitting occurs why?

→ Because though it has low bias but it will have high variance. new test data is added to avoid this learning rate is introduced ( $\alpha = 0-1$ )

Still 73.7 is quite large than  $y$  so next decision tree will be added based on this error.

$$F(x) = h_0(x) + \alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots = \sum_{i=1}^n \alpha_i h_i(x)$$



Machine Learning is nothing but finding relationship between %p and output data.

### Additive Modelling

Adding functions/DT at every step and boosting Algorithm uses additive modelling format.

### Pseudo Algorithm Explanation of Gradient Boosting

Ex:-

<u>R&amp;D Spend</u>	<u>Administration</u>	<u>Marketing Spend</u>	<u>(Y) profit</u>
165	137	472	192
101	92	250	144
29	127	201	91

Step 1:- Initialize

#### Requirements

- 1)  $x_i, y_i$   
 $\downarrow$   
 independent features  $\rightarrow$  dependent features
- 2) differentiable Loss function  $L(y, f(x))$  based on either Regression or classification problem.

$$L(y, f(x)) \rightarrow L(y, \hat{y})$$

$\uparrow$   
 $\hat{y}$

Here we are using MSE (Mean Square Error)  $= \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

Step 1:- Initialize  $f_0(x) = \arg \min_x \sum_{i=1}^n L(y_i, x)$

**Aim:-** our aim is to establish relationship between  $y = f(x)$   
 Since, it is a boosting  $F(x)$  is obtained adding small functions sequentially.

$$F(x) = f_0(x) + f_1(x) + f_2(x) + \dots + f_n(x)$$

Base model

$\sum$   
DT

Step 1 :- Initialize  $F_0(r) = \arg \min_r \sum_{i=1}^n L(y_i, r)$  / Initialize Base Model.

$$r = \hat{y}$$

$$\text{Loss function} = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y})^2$$

$$F_0(r) = \arg \min_r \frac{1}{2} \sum_{i=1}^n (y_i - r)^2$$

Here our aim is to get "r" value in such a way that the error is less.

↓ diff w.r.t r

$$\frac{d F_0(r)}{dr} = \frac{d}{dr} \frac{1}{2} \sum_{i=1}^n (y_i - r)^2 = \frac{1}{2} \sum_{i=1}^n \frac{d}{dr} (y_i - r)^2$$

$$= \sum_{i=1}^n (y_i - r) \frac{d}{dr} (y_i - r)$$

$$0 = - \sum_{i=1}^n (y_i - r)$$

$$\sum_{i=1}^n (r - y_i) = 0$$

$$\sum_{i=1}^3 (r - y_i) = 0$$

$$(r - 192) + (r - 144) + (r - 91) = 0$$

$$3r = 192 + 144 + 91$$

$$r = \frac{192 + 144 + 91}{3} \quad ({}^\circ \circ \text{ mean of the o/p column})$$

$$= \frac{427}{3}$$

$$= 142.33$$

Here, the base model value is 142.33



② For  $m = 1$  to  $M$ :

$$F(x) = \underbrace{f_0(x)}_{\substack{\downarrow \\ \text{just a} \\ \text{leaf} \\ \text{Base model}}} + \underbrace{f_1(x)}_{DT} + \underbrace{f_2(x)}_{DT} + \dots + f_M(x)$$

$M$  = no. of Decision Tree

a) for  $i = 1, 2, \dots, N$  compute

$$r_{im} = - \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}$$

$i$  = no. of rows

$m$  = on which DT

we are working

$r$  = residual / Pseudo Residual

**Aim**:- need to calculate Pseudo Residual for every Decision Tree.

let  $m=1$

$$r_{i1} = - \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_0}$$

You must have noticed in this Dataset we have three rows  
So, we need to calculate:-

$r_{11}, r_{21}, r_{31}$   $\rightarrow$  1st DT, 2 row Pseudo Residual  
 $\downarrow$   
1st DT, 1 row Pseudo Residual

$$f(x_i) = \hat{y}_i$$

$$\text{Loss function} = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$r_{i1} = - \left[ \frac{\partial L(y_i, \hat{y}_i)}{\partial \hat{y}_i} \right]_{f=f_0}$$

$$= - \left[ \frac{\partial}{\partial \hat{y}_i} \frac{1}{2} (y_i - \hat{y}_i)^2 \right]_{f=f_0}$$

$$= - \left[ \frac{1}{2} \times 2 (y_i - \hat{y}_i) \frac{\partial}{\partial \hat{y}_i} (y_i - \hat{y}_i) \right]_{f=f_0}$$

$$= - [(y_i - \hat{y}_i) (0 - 1)]_{f=f_0}$$

$$= [y_i - \hat{y}_i]_{f=f_0}$$

= observed - predicted

$$= [y_i - f(x_i)]_{f=f_0} = [y_i - f_0(x_i)]$$

$$r_{11} = y_1 - f_0(x_1) = 192 - 142 = 50$$

$$r_{21} = y_2 - f_0(x_2) = 144 - 142 = 2$$

$$r_{31} = y_3 - f_0(x_3) = 91 - 142 = -51$$

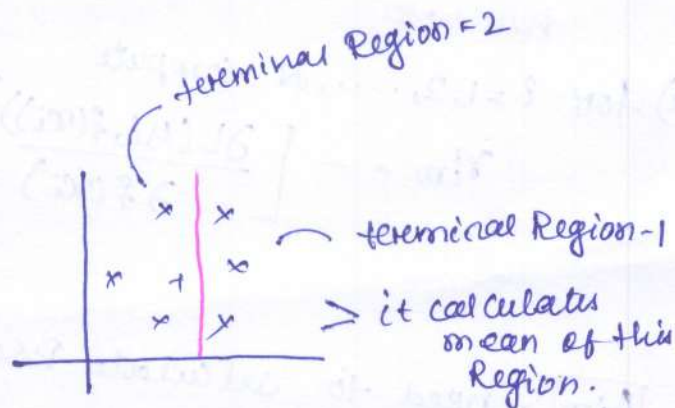
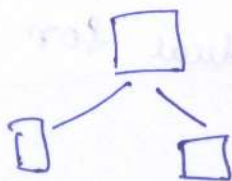
2(b) fit a regression tree to the targets  $Y_{im}$  giving terminal regions  $R_{jm}, j=1, 2, \dots, J_m$

i.e

Considering three independent columns as ~~the~~ input and  $Y_{im}$  as o/p we need to fit a tree.

Here, what is terminal Region?

Suppose we have  $X|Y$  and DT is like this



It will cut Graph in to two halves as the DT i.e one depth we have taken

$R_{jm} = \text{terminal Region}$

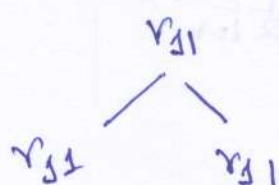
2(c) for  $j=1, 2, \dots, J_m$ , compute

$$r_{jm} = \arg \min_r \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + r)$$

of terminal Regions

Here, the aim is to calculate the o/p based on above formula

$$r_{j1} = \arg \min_r \sum_{x_i \in R_{j1}} L(y_i, f_{m-1}(x_i) + r)$$



$x_i \in R_{j1} \Rightarrow R_{j1}$  is terminal region

— This means we need to calculate the points that are fallen into that terminal region only.

$$r_{j1} = \arg \min_r \frac{1}{2} (y_i - (f_0(x_i) + r))^2$$

$$\begin{aligned} \frac{dL}{dr} &= \frac{1}{2} \times 2 \times (y_i - f_0(x_i) + r) \frac{d}{dr} (y_i - f_0(x_i) + r) = 0 \\ &= (y_i - (f_0(x) + r)) (\pm 1) \end{aligned}$$

$$\begin{aligned} r_{j1} &= y_i - f_0(x) - r \\ &= 91 - 142 - r = 0 \\ r &= -51.00 \end{aligned}$$



# XG-Boost/Extreme Gradient Boosting

(13)

> XGBoost improves upon the base of GBM framework through system optimization & algorithmic enhancements.

## System optimization

- 1) Parallelization.
- 2) Tree pruning
- 3) Hardware optimization

## Algorithmic Enhancements

- 1) Regularization
- 2) Sparsity awareness
- 3) Weighted Quantile Sketch
- 4) Cross validation.

## XGBoost classifier

Salary	Credit	Approval	Res	New pr
<=50K	Bad	0	0 - Pr = -0.5	0.6
<=50K	Good	1	1 - 0.5 = 0.5	-
<=50K	Good	1	0.5	-
>50K	Bad	0	-0.5	-
>50K	Good	1	0.5	-
>50K	Normal	1	0.5	-
<=50K	Normal	0	-0.5	-

ep 1:- Construct a Base model

o/p will be 0.1 =  $\frac{0+1}{2} = 0.5 = P_r$

Construct a tree

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

Salary

<=50K >50K

S.W = 0.33

S.W = 0  
 $[-0.5, 0.5, 0.5, -0.5]$

$[-0.5, 0.5, 0.5]$

$$\frac{0.5 + 0.5 + 0.5 - 0.5}{5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5)} = 0$$

$$S.W = \frac{[-0.5 + 0.5 + 0.5]^2}{(0.5(1-0.5)) \times 3} = \frac{0.25}{(0.25) \times (0.25) \times (0.25)} = 0.33$$

Calculate the Similarity weight (S.W)

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

Calculate Gain = (S.W of left hand side + S.W of RHS) - S.W of Root node

$$= (0 + 0.33) - 0.1428 = 0.1872$$

Which ever has more gain takes that split.  
 Suppose Gain is more for Salary so it takes the first split.

Here also we need to check which split it is taking and as per the gain we do the split.

$$\text{Gain} = (1 + 0.33) - 0 = 1.33$$

$$S.W = \frac{0.25}{0.5 \times 0.5} = 1$$

$$S.W = \frac{(0.5)^2}{(0.25) \times 3} = 0.33$$

$$\text{Gain} = 1$$

→ Calculate the o/p of Base Model.

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

$$= 0$$

$$= \sigma\left(0 + \overset{0.1}{L_1} (\text{similarity rate of } DT_1)\right) = \sigma(0 + 0.1 \times 1) = \sigma(0.1)$$

Basemodel

→ Sigmoid Activation function =  $\frac{1}{1+e^{-x}} = \frac{1}{1+e^{-0.5}} = 0.6$

$$XGBOOST = \underbrace{\text{Base model}}_{\log\left(\frac{P}{1-P}\right)} + \sigma(L_1(DT_1) + L_2(DT_2) + \dots)$$



# XG-Boost Regressor

(15)

Ex:-

<u>Exp</u>	<u>Grp</u>	<u>Salary</u>	<u>Res</u>
2	Yes	40K	$40 - 51.2 = -11$
2.5	Yes	42K	-10
3	No	52K	1
4	No	60K	9
4.5	Yes	62K	11

① Create the Base model

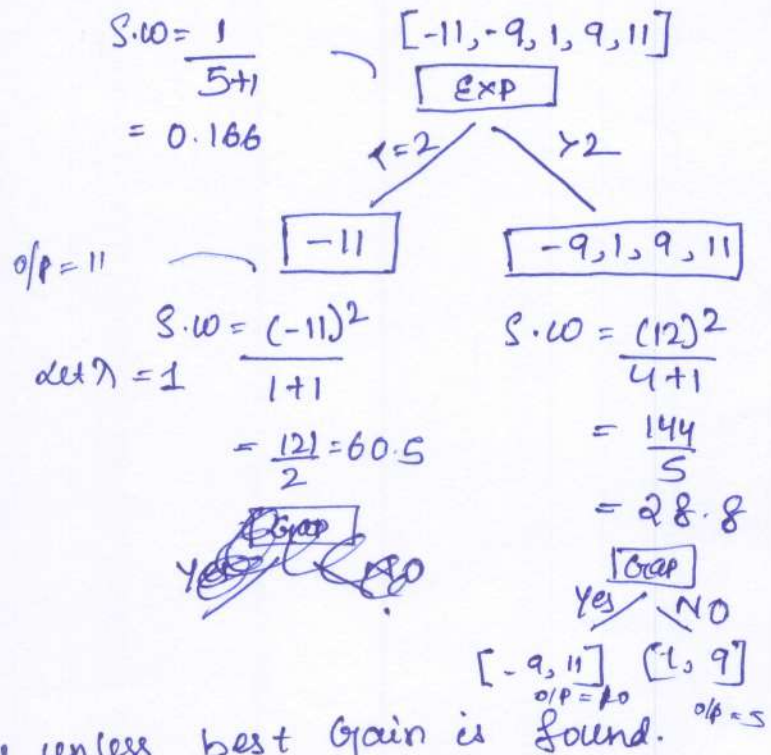
$$= \frac{40 + 42 + 52 + 60 + 62}{5}$$

$$= 51.2K$$

$$S.W = \frac{\sum (\text{Residual})^2}{\text{No. of Residual} + 1}$$

$$\text{Gain} = (60.5 + 28.8) - 0.166$$

$$= 89.134$$



∴ like way splitting happens until unless best Gain is found.

$$o/p = \alpha_1 T_1 + \alpha_2 T_2 + \alpha_3 T_3$$

↓                      ↓

Base model          avg of Decision Tree.