

lets modify the sample weights so that next stump will take into account the error made by first stump.

To Modify weight

$$\text{New Sample weight} = \text{Sample weight} \times e^{-\text{amount of say}}$$

of incorrectly classified sample

Note

if Amt of say is less than the stump is not doing good job -

In our case for incorrectly classified sample

$$\text{New Sample weight} = \frac{1}{8} \times e^{0.97} = 0.33$$

for correctly classified sample

$$\text{New Sample weight} = \text{Sample weight} \times e^{-\text{amount of say}}$$

$$\text{Updated Sample weight} = \frac{1}{8} \times e^{0.97} \Rightarrow 0.05$$



Now

Chest Pain	Blocked Arteries	Patient weight	Heart Disease	Sample weight	New Sample weight	Normalized Sample weight
Yes	Yes ✓	205	Yes	1/8	0.05	0.07
No.	Yes ✓	180	Yes	1/8	0.05	0.07
Yes-	No -	210	Yes	1/8	0.05	0.07
Yes-	Yes ✓	167	Yes	1/8	0.33	0.49
No.	Yes -	156	No	1/8	0.05	0.07
No.	Yes -	125	No	1/8	0.05	0.07
Yes-	No	168	No	1/8	0.05	0.07
Yes-	Yes -	172	No	1/8	0.05	0.07

stump 2 will force on this more.

Yes -	No -	$\frac{1}{17}$	No -	$\frac{1}{8}$	0.05	
Yes -	No -	$\frac{1}{17}$	No -	$\frac{1}{8}$	0.05	

Normalized Not
Normalized Total sum = 0.68
Divide Group
Entire in N=10
Sample weight by 0.68

We have new Normalized Sample weight \rightarrow

The ~~incorrectly~~ classified Sample has Higher New Sample weight

\rightarrow which means Stump2 will focus more on it.

\rightarrow The correctly classified sample has comparatively lower sample weight so Stump2 will have less focus on it

* Similarly we can have many stumps.

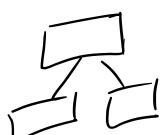
. This is how the Adaboost creates & uses stumps.

Now lets chk how forest of stumps made by Adaboost does classification.

(1) We will pass the test sample through each stump.

(2) Let there be few stumps that classifies the sample as Heart Disease & let there be stumps that classifies the sample as Not Heart Disease

(3) Has Heart Disease

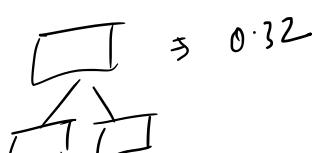


\Rightarrow Amt of Say
0.97

Has Not Heart Disease



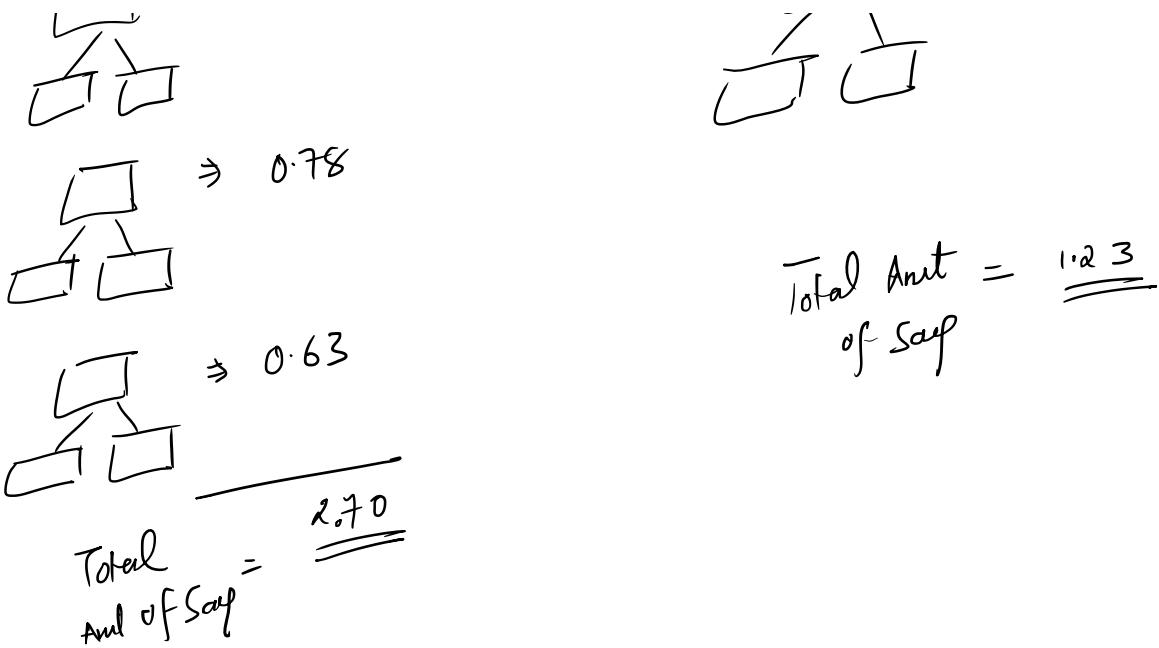
\Rightarrow Amt of Say
0.41



\Rightarrow 0.32



\Rightarrow 0.82



$\overline{\text{Total Amt of Samp}}$ For Stumps that classifies as Heart Disease is greater
 \therefore The test sample is classified as Heart Disease

Summary

- ① Assign Sample Weights (Initial)
- ② Create Stumps for each Property
- ③ Use gini index to identify first stump
- ④ For the first stump
 - calculate Total Error
 - calculate Amt of Samp
- ⑤ update the sample weights for each sample
 - For Incorrectly Classified = $\frac{\text{Updated weight}}{\text{weight}} = \text{Sample weight} \times e^{\text{amt of samp}}$
 - For Correctly Classified = $\frac{\text{Updated weight}}{\text{weight}} = \text{Sample weight} \times e^{-\text{amt of samp}}$
- ⑥ Now Specify Updated weight for each sample & Normalize it

to generate New Sample Weight

(7) Identify New Stump based on New Sample weight
& Repeat step 3 to 7

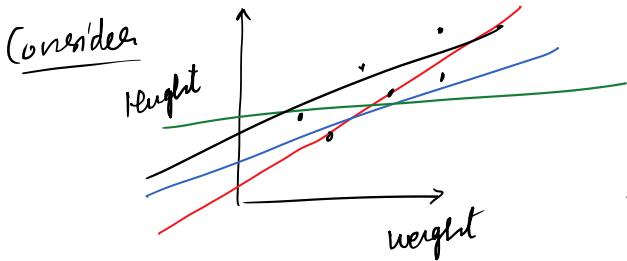
(8) For Classification

8.1 Run the Test Sample through All the Stumps

8.2 Calculate Total Amt of Sump for Samples classifying Yes

8.3 & NO Classify the test sample based on largest sum of Amt of Sump.

Gradient Descent →



$$\text{predicted} = \text{intercept} + \text{slope} * \text{weight}$$

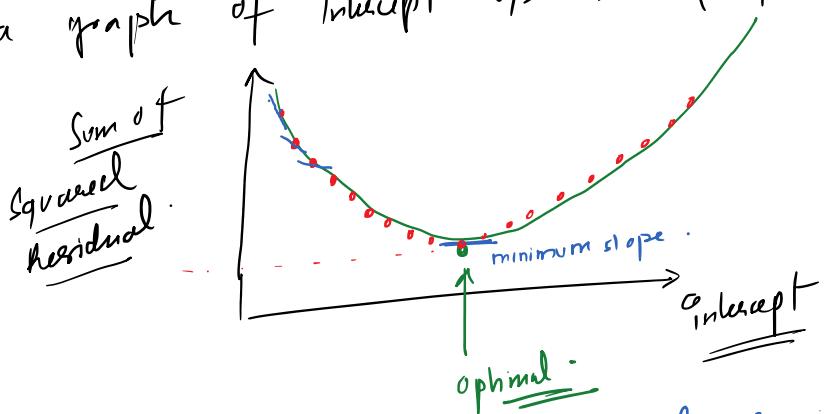
Here we can Many line with diff
intercept & slope values.

We want line with intercept & slope such that the cost J^n should give minimum value.

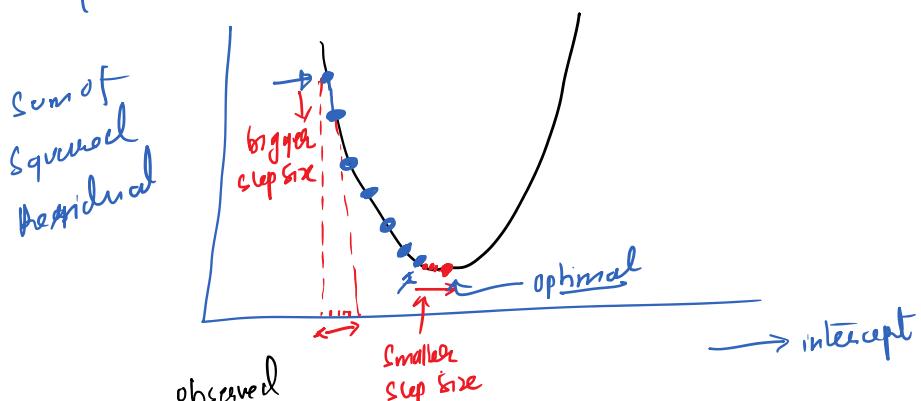
Let the cost J^n be sum of squared error (Residual)

Let us consider only intercept

Let's plot a graph of Intercept v/s Sum of Squared Residual



To find Minimum value for sum of squared residual we will start with small value of intercept and will increase the value by very small amount.



observed

Smaller Step size

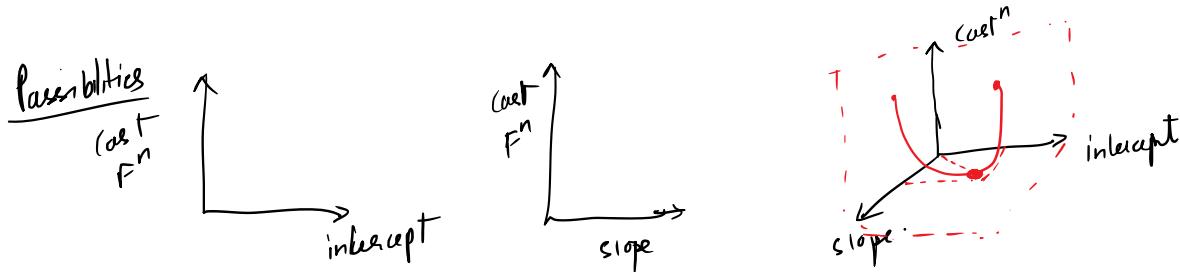
→ intercept

Consider

x	y
0.5	1.4
2.3	1.9
2.9	3.2

Let the cost $F^n = \frac{\text{Sum of Squared Residual}}{\text{observed}}$

$$\begin{aligned} \text{Sum of Squared Residual} &= \left(\frac{1.4 - (\text{intercept} + \text{slope} * 0.5)}{\text{observed}} \right)^2 \\ &\quad + \left(1.9 - (\text{intercept} + \text{slope} * 2.3) \right)^2 \\ &\quad + \left(3.2 - (\text{intercept} + \text{slope} * 2.9) \right)^2 \end{aligned}$$



Here we want to find value of f intercept & slope that gives Minimum Sum of Squared Error.

Slope $\frac{d}{d \text{slope}} (\text{Sum of Squared Residual}) = \frac{d}{d \text{slope}} (\text{Sum of Squared Residual}).$

Intercept

Using chain rule.

$$\frac{d}{d \text{intercept}} \left(\frac{1.4 - (\text{intercept} + \text{slope} * 0.5)}{\text{slope}} \right)^2$$

$$= 2 \times \frac{1.4 - (\text{intercept} + \text{slope} * 0.5)}{\text{slope}} * \frac{d}{d \text{intercept}} (\text{intercept} + \text{slope} * 0.5)$$

$$= -2 \left(1.4 - (\text{intercept} + \text{slope} * 0.5) \right).$$

$$\frac{d}{d \text{intercept}} (\text{Sum of Squared Residual}) = (-2)(1.4 - (\text{intercept} + \text{slope} \times 0.5)) \\ + (-2)(1.9 - (\text{intercept} + \text{slope} \times 2.3)) \\ + (-2)(3.2 - (\text{intercept} + \text{slope} \times 2.9)). \quad \left. \right\} \text{Eq 1}$$

$$\frac{d}{dslope} (\text{Sum of Squared Residual}) = (-2)(0.5)(1.4 - (\text{intercept} + \text{slope} \times 0.5)) \\ + (-2)(2.3)(1.9 - (\text{intercept} + \text{slope} \times 2.3)) \\ + (-2)(2.9)(3.2 - (\text{intercept} + \text{slope} \times 2.9)). \quad \left. \right\} \text{Eq 2}$$

We got the above 2 Eqⁿ let's start with random values of
 $\text{intercept} = 0$ & $\underline{\text{slope}} = 1$

Substitute the values $\rightarrow \text{intercept} = 0 \quad \underline{\text{slope}} = 1$

$$\frac{d}{d \text{intercept}} = -1.6 \quad \left. \right\} \text{slope of line w/o intercept}$$

$$\frac{d}{dslope} = -0.8 \quad \left. \right\} \text{slope of line w/o slope}$$

for step size
 New stepsize = slope of line at point \times learning Rate

Note: learning Rate \Rightarrow takes small value $0.1 \rightarrow \underline{0.3}$

Let us take learning Rate = 0.01

$$\therefore \text{stepsize}_{\text{intercept}} = \frac{d}{d \text{intercept}} \times 0.01 = -1.6 \times 0.01 = \underline{-0.016}$$

$$\text{stepsize}_{\text{slope}} = \frac{d}{dslope} \times 0.01 = -0.8 \times 0.01 \Rightarrow \underline{-0.008}$$

$$\text{step size} = \frac{\alpha}{\text{slope}} * \text{Old Val}$$

$$\text{New Intercept} = \text{Old Intercept} - \frac{\text{step size}}{\text{slope}} = 0 - (-0.016) = \underline{\underline{0.016}}$$

$$\text{New slope} = \text{Old slope} - \frac{\text{step size}}{\text{slope}} = 1 - (-0.008) = \underline{\underline{1.008}}$$

Substitute this New Intercept & New slope in Eqⁿ $\frac{d}{d\text{Intercept}}, \frac{d}{d\text{Slope}}$
 & calculate new step size for intercept & slope
 And again calculate New Intercept & New Slope

We will repeat until step size becomes very small = $\frac{0.001}{\epsilon_x}$
 or some maximum number of steps is reached i.e. $\underline{\underline{1000}}$

for above Dataset Best fitting line will have Intercept = 0.95
 & slope = 0.64

This is how the gradient descent optimizes parameters.

Summary

- ① Get Eqⁿ of line
- ② Want best Intercept & slope
- ③ Identify loss fn (Sum of Squared Residual)
- ④ Find Eqⁿ for $\frac{d}{d\text{Intercept}}$ & $\frac{d}{d\text{Slope}}$
- ⑤ Start with initial value of slope & intercept
 ~ 0.01 & $\frac{d}{d\text{...}}$

(5) Start with

(6) find $\frac{d}{d \text{intercept}}$ & $\frac{d}{d \text{slope}}$

$$(7) \text{ calculate } \frac{\text{stepSize}_{\text{intercept}}}{\text{intercept}} = \frac{d}{d \text{intercept}} * \text{learningRate}$$

$$\frac{\text{stepSize}_{\text{slope}}}{\text{slope}} = \frac{d}{d \text{slope}} * \text{learningRate}$$

$$(8) \text{ Updated Intercept} = \text{old Intercept} - \text{stepSize}_{\text{intercept}}$$

$$\text{Updated Slope} = \text{old Slope} - \text{stepSize}_{\text{slope}}$$

(9) repeat step 6 \rightarrow step 8 until step size is very small

or No of iterations $\geq \underline{\underline{1000}}$

Gradient Boost for Regression \rightarrow

- * Gradient Boost starts with single leaf instead of Tree/stump.
- * Leaf represents initial guess for the weights of the samples - observation

Ex Consider

Height	favorite color	gender	<u>Weight</u>
1.6	Blue	male	88
1.6	Green	female	76
1.5	Blue	Female	56
1.8	Red	male	73
1.5	Green	male	77
1.4	Blue	female	57

Initial Prediction

Step 1 First we make Initial Guess = weight (Average) = $\underline{71.2}$

Step 2 Like Adaboost, even Gradient Boost builds tree based on previous tree

But Unlike Adaboost, the gradient Boost tree are larger than stump

However the Height is still restricted [No of leaf $8 \rightarrow 32$ in each Tree]

[Gradient Boost will build another tree based on error of previous tree]

- * Gradient Boost continues to build tree in this fashion until it has made numbers of trees you asked for or the Additional tree fails to improve the fit

or the Additional tree turns to imp.....

Here Initial Guess = weight (Average) = 71.2

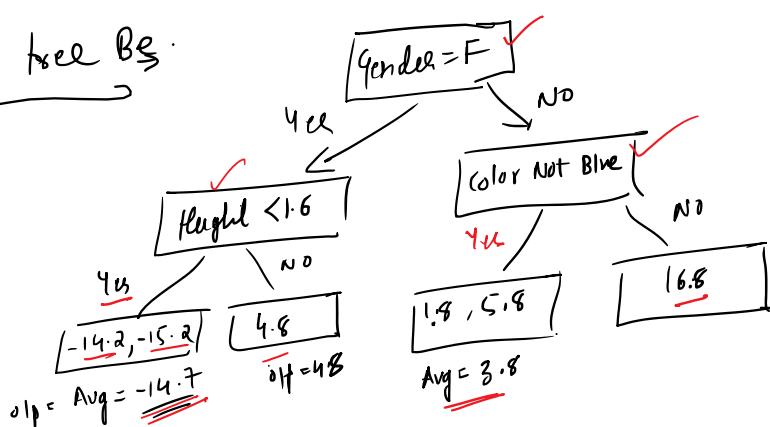
$$\text{Predicted Value} = \boxed{71.2} + \underset{\substack{\uparrow \\ \text{Lat}}}{\text{---}} + \underset{\substack{\longrightarrow \\ \text{Boosting}}}{\text{---}}$$

Now we need to find residual for all the samples.

Height	favorite color	Gender	Weight	Residual
1.6	Blue	Male	88	16.8
1.6	Green	Female	76	-4.8
1.5	Blue	Female	56	-15.2
1.8	Red	Male	73	1.8
1.5	Green	Male	77	5.8
1.4	Blue	Female	57	-14.2

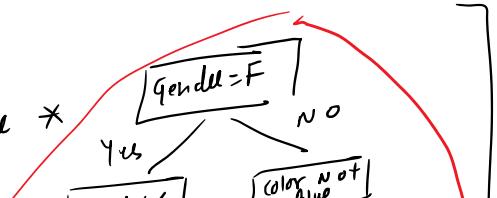
Let's build a Tree using the Residual & not weight

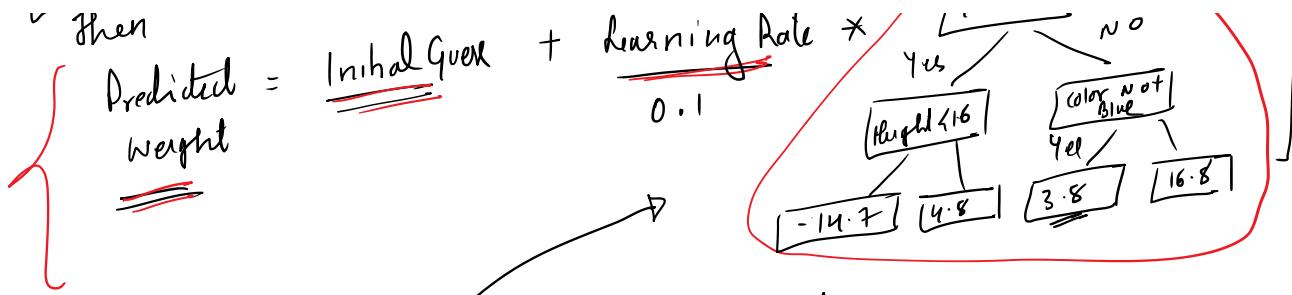
Let the tree Be:



Predict If there will be only 1 Tree Assume

$$\text{Predicted} = \text{Initial Guess} + \text{Learning Rate} * \frac{1}{n} \sum_{i=1}^n \text{Residual}_i$$





Ex for first sample

1.6	Blue / Male	2
-----	-------------	---

$$\text{Predicted weight} = 71.2 + 0.1 * 3.8 = 71.58$$

Apply this on All the samples to get New Prediction

Height	favorite color	Gender	Weight	Residual	New Prediction	New Residual
1.6	Blue	Male	88	16.8	71.58	R ₁
1.6	Green	Female	76	4.8	P ₁	R ₂
1.5	Blue	Female	56	-15.2	P ₂	R ₃
1.8	Red	Male	73	1.8	P ₃	R ₄
1.5	Green	Male	77	5.8	P ₄	R ₅
1.4	Blue	Female	57	-14.2	P ₅	R ₆

Residual 1

Residual 2

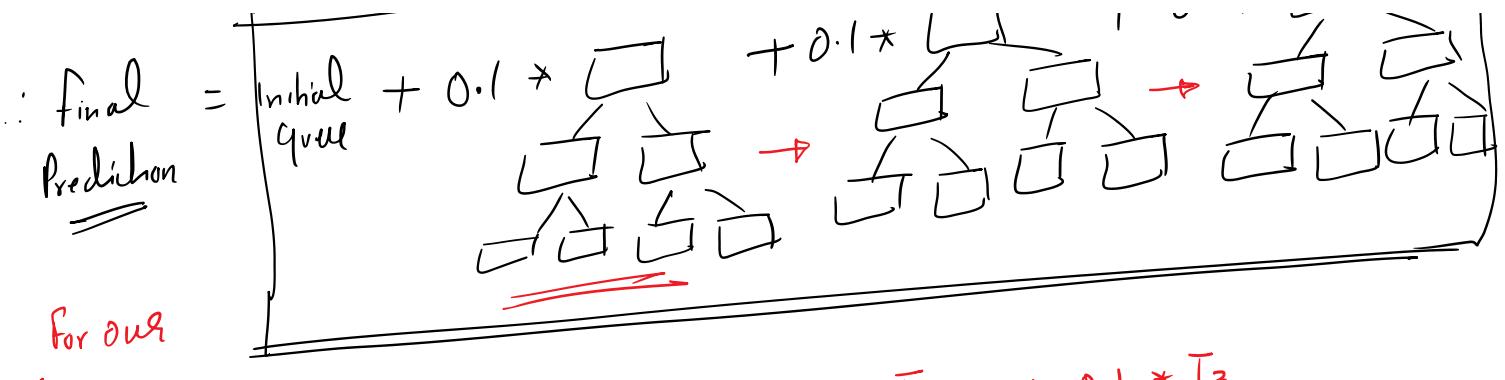
* Now Again Find New Residual for New Prediction for all the samples -

* Using this New Residual Construct New Tree and So on.

Until the New Tree does not make significant contribution to New Prediction.

Let say we have 3 Tree for above dataset

$$\text{Final} = \text{Initial} + 0.1 * \boxed{\quad} + 0.1 * \boxed{\quad} + 0.1 * \boxed{\quad}$$



for our ex

$$= 71.2 + 0.1 * T_1 + 0.1 * T_2 + 0.1 * T_3$$

Boosting

Gradient Boost Principle →

[Taking lot of small steps in right direction will result into better predictions]

* Gradient Boost for Classification →

* Consider

Observation (target)			
likes Popcorn	Age	Favorite color	loves <u>Troll 2</u> (movie)
Yes	12	Blue	Yes 1
Yes	87	Green	Yes 1
No	44	Blue	No 0
Yes	19	Red	No 0
No	32	Green	Yes 1
No	14	Blue	Yes 1

C Classification

Step 1 → Start with leaf that represent initial prediction.

* When we use Gradient Boost for classification the initial prediction for every individual sample = log(odd) [In regression, it's average of the observation]

log(odd) → That someone loves Troll 2 movie

Here 4 ppl loves Troll 2 & 2 do not love Troll 2

$$\log_e(\text{odd}) = \log\left(\frac{\text{no of Yes}}{\text{no of No}}\right) = \log_e\left(\frac{4}{2}\right) = \ln(2) = \boxed{0.7}$$

To use log(odd) we have to convert it into probability.

We will use logit function

$$\therefore \text{Probability of loving Troll 2} = \frac{e^{\log(\text{odd})}}{1 + e^{\log(\text{odd})}} = \frac{e^{0.7}}{1 + e^{0.7}} = \frac{1}{1 + e^{-0.7}} = \boxed{0.7}$$

$$\text{Probability} = \frac{\log(\text{odd})}{1+e^{-x}} = \frac{1+e^x}{1+e^{x+\log(\text{odd})}}$$

[Since probability of loving Trolls is greater than 0.5, so we can classify everyone in Training Dataset as someone who loves Trolls]

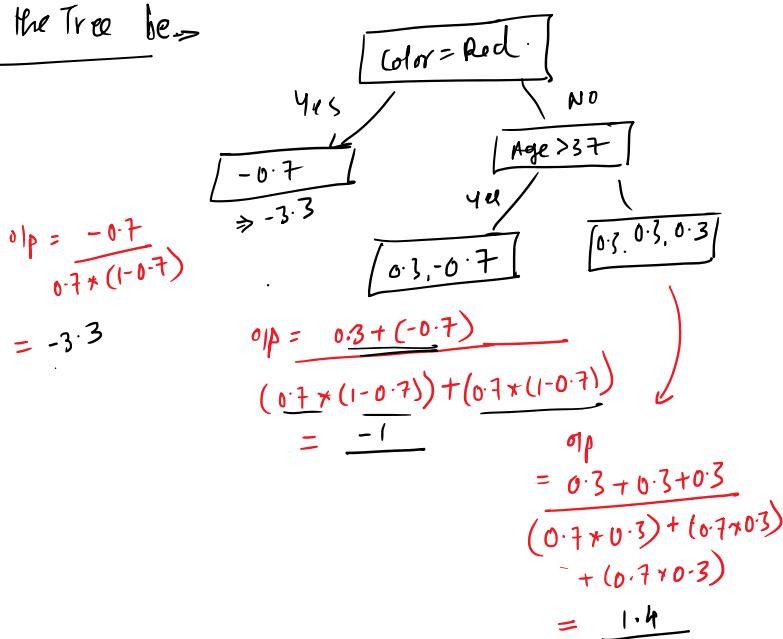
Likes Popcorn	Age	Favourite Color	Hates Trolls (movie)	Residual
Yes	12	Blue	Yes 1 ✓	0.3
Yes	87	Green	Yes 1	0.3
No	44	Blue	No 0	-0.7
Yes	19	Red	No 0	-0.7
No	32	Green	Yes 1	0.3
No	14	Blue	Yes 1	0.3

Initial = 0.7

Initial Residual

Now Build Tree Using Color and Age.

Let the Tree be →



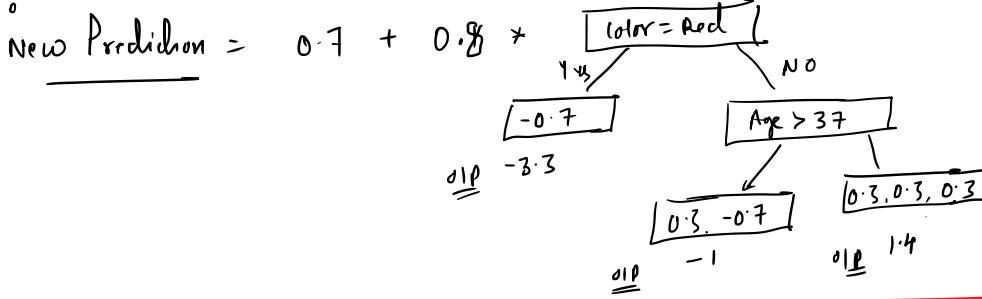
Note → When gradient boost is used for regression the leaf with single residual has OLP equal to that residual & for leaf with more than one residual the OLP is average of residual.

For classification:

OLP value at each leaf

= $\sum \text{Residual}$

$\sum (\text{Previous Prob} \times (1 - \text{Previous Prob}))$
For each residual.



Consider first sample $\begin{bmatrix} \text{Likes} & \text{Age} & \text{Color} \\ \text{Yes} & 12 & \text{Blue} \end{bmatrix}$.

$$\text{New Prediction} = 0.7 + 0.8 \times 1.4 = \underline{\underline{1.8}}$$

$$\text{Convert into Probability} = \frac{1}{1+e^{-1.8}} = \underline{\underline{0.9}}$$

Now find New Prediction for each sample using above model.

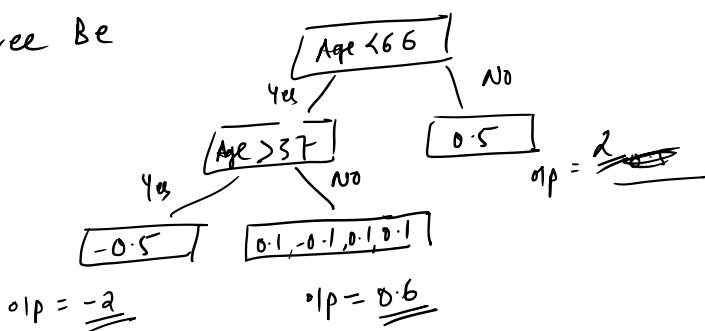
✓ observation Initial ~~Prediction = 0.7~~ let

Likes Popcorn	Age	Favourite color	Observation	Initial Prediction = 0.7	New Prediction	New Residual
Yes	12	Blue	Yes 1 ✓	0.3	0.9 ✓	0.1
Yes	87	Green	Yes 1	0.3	0.5	0.5
No	44	Blue	No 0	-0.7	0.5	-0.5
Yes	19	Red	No 0	-0.7	0.1	-0.1
No	32	Green	Yes 1	0.3	0.9	0.1
No	14	Blue	Yes 1	0.3	0.9	0.1

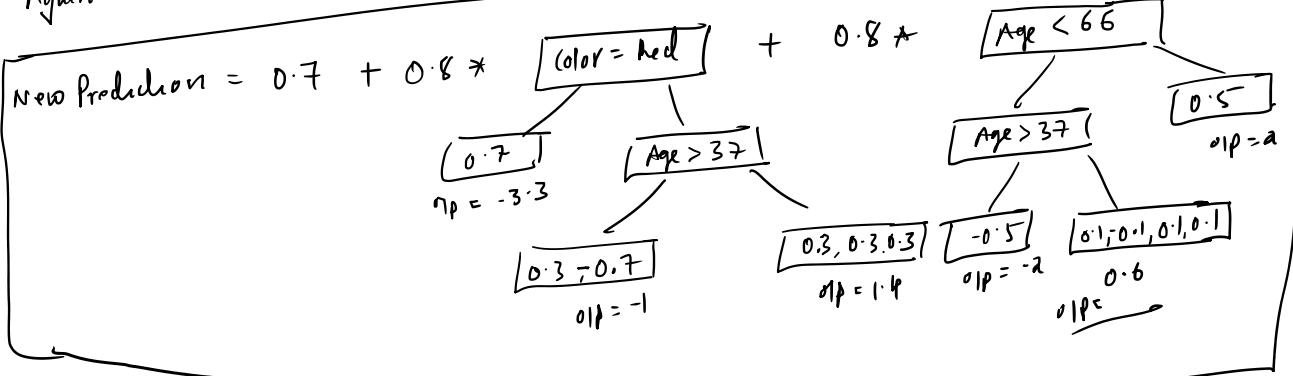
Here New Prediction is not 0.7 for all sample, then each sample has its own new prediction & we calculated New Residual for each Sample.

Again Build a Tree Using ut Age only

Let the Tree Be



Again find New Probability Prediction.



Again apply each Sample detail on the above model and calculate New prediction.

→ Again calculate new residual and so on.

This process repeats until we have made the maximum specified no of Tree or the residual gets very very small

Summary Gradient Boost for Classification

① Consider the dataset

② Find $\log(\text{odd})$.

③ Find Initial Probability using $\log(\text{odd}) = \frac{1}{1 + e^{-\log(\text{odd})}}$

④ Calculate Residual based on Initial Probability

⑤ Construct Tree using Initial Residual.

⑥ New Prediction = Initial Prediction + LearningRate * Tree.

⑦ Calculate New Prediction for all Sample.

⑧ Again find New Residual.

⑨ Construct New Tree and so on.

⑩ Repeat until max no of Tree constructed or the Residual is Insignificant.

Summary Gradient Boost for Regression

- ① Consider the dataset
- ② Find Initial prediction = Avg of the Observed value.
- ③ Calculate Residual based on Initial Prediction
- ④ Construct Tree Using Initial Residual.
- ⑤ New Prediction = Initial Prediction + LearningRate * Tree.
- ⑥ Calculate New Prediction for all Sample.
- ⑦ Again find New Residual.
- ⑧ Construct New Tree and so on.
- ⑨ Repeat until max no of Tree constructed or the residual is Insignificant.

XGBoost (Extreme Gradient Boost)

→ Implementation of Gradient Boosted Decision Tree
Here Decision Tree is created in Sequential form.

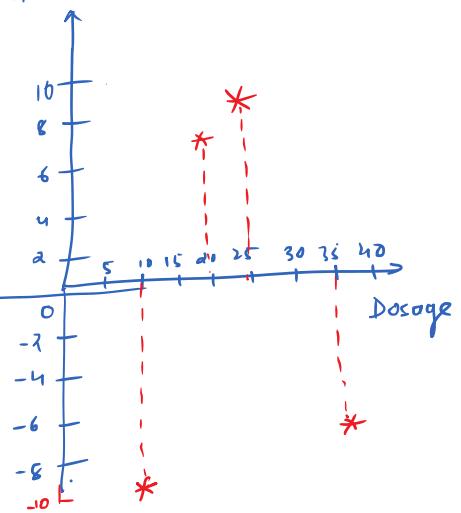
- Weights plays important role in XGBoost
- Weights are assigned to all independent variables which are then fed into Decision Tree that predict result.
- Weights of variable predicted wrong by tree is increased and that variable are fed in second decision tree.
- These individual predictors/classifiers can ensemble to give strong and more precise model.
- XG Boost is faster than Gradient Boost
- There is stopping criterion for tree splitting in XG Boost.
 - XG Boost uses max_depth parameter and it starts pruning the tree backward.
- This pruning improves computational performance and also helps to overcome problem of overfitting

XGBoost for Regression

Consider

x (Dosage of Drug)	y (Drug Effectiveness)
10	-10
20	7
25	8
35	-7

Effectiveness



Build Prediction Model Using XGBoost.

Step 1 → Assume ^{initial} threshold for effectiveness = 0.5^-

Step 2: Calculate Residual.

x	y	observation	Initial = 0.5	Residual	Initial Prediction
10	-10			-10.5	
20	7			6.5	
25	8			7.5	
35	-7			-7.5	

Initial Residual

Step 3 Consider all the residuals in knot, let $\lambda = 0$

$$\boxed{-10.5, 6.5, 7.5, -7.5}$$

$$\text{Similarity} = \frac{(-10.5 + 6.5)^2}{4+0} = 4$$

Now let us see if we can cluster the residual better on Similarity by split

Note

$$\text{Score} = \frac{(\text{Sum of Residual})^2}{\text{No of Residual} + \lambda}$$

λ = Regularization Parameter

$$\text{Gain} = \text{Similarity}(\text{left}) + \text{Similarity}(\text{right}) - \text{Similarity}(\text{root})$$

To Determine Dosage value for Split \Rightarrow

Case 1) Consider Avg of first two x

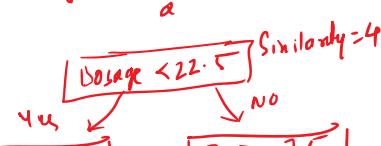
$$\text{Average Dosage} = 15$$



Case 2)

Consider Avg of Next two x

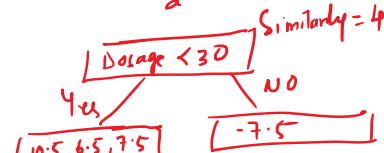
$$\text{Avg} = \frac{20+25}{2} = 22.5$$

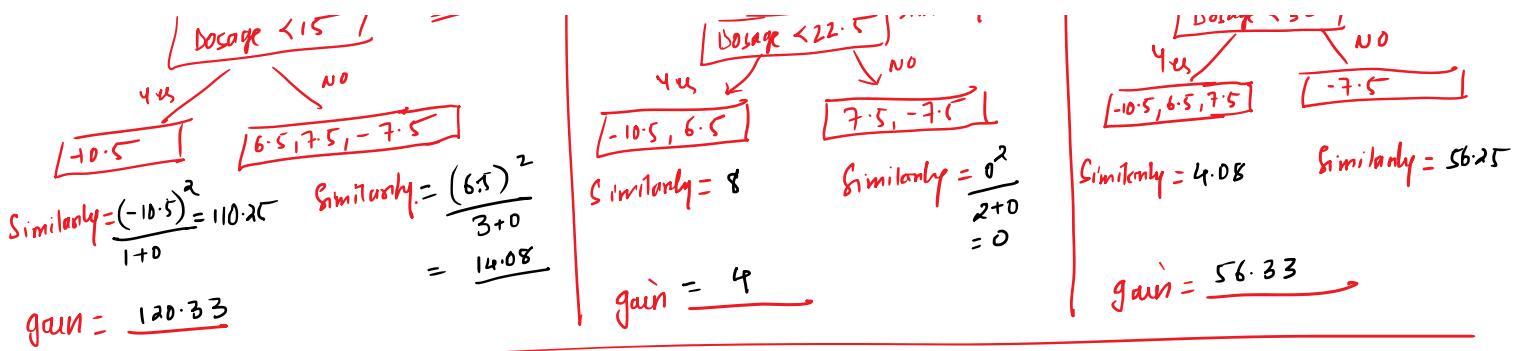


Case 3)

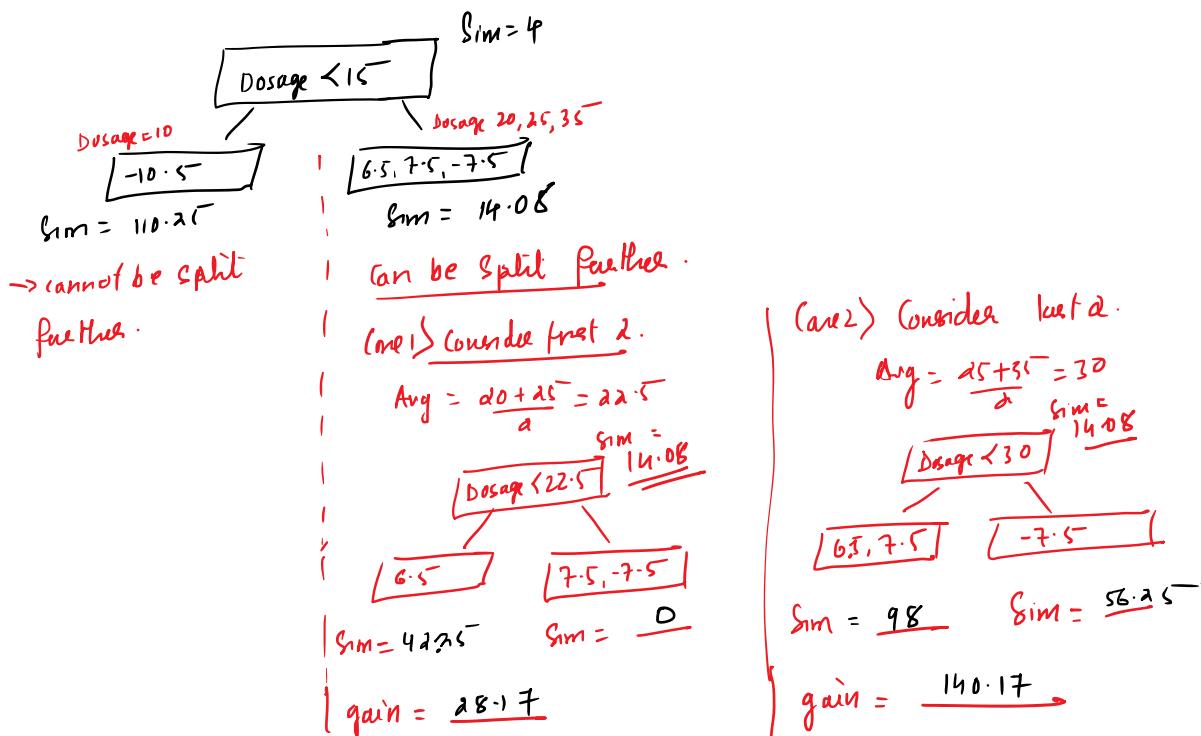
Consider Avg of last two x

$$\text{Avg} = \frac{25+35}{2} = 30$$

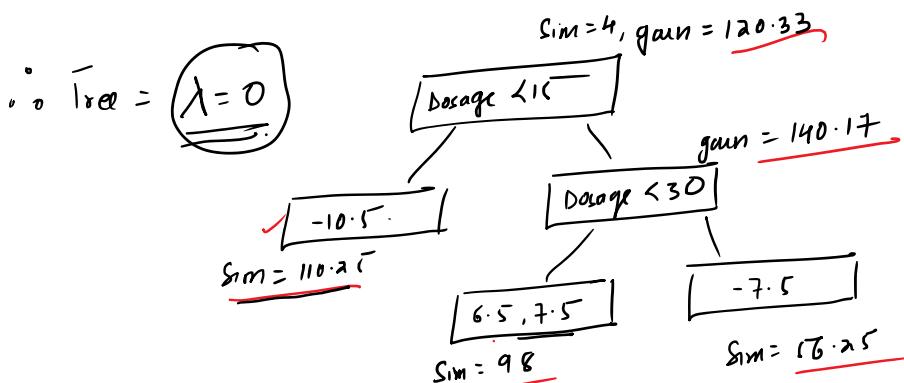




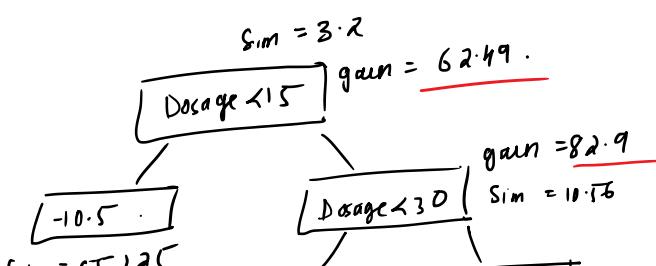
Since $\text{Dosage} < 15$ gives highest gain so it will be split factor.

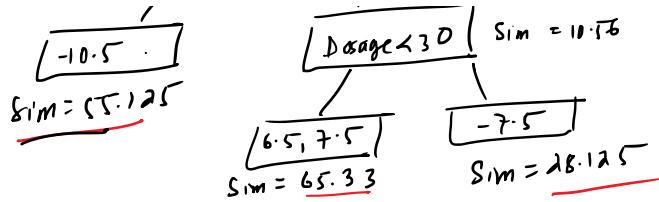


Highest gain is for $\text{Dosage Avg} = 30$. So
 $\text{Dosage} = 30$ will be split condition.



Solve it for $\underline{\lambda = 1}$ Tree





Note
when $\lambda > 0$ the Similarity Score is smaller so gain is smaller.

[Decrease in Similarity Score is inversely proportional to Number of individuals in the leaf].

Pruning \rightarrow Kyon chahiye.

- ① Height restriction for Tree
- ② Ignoring Unnecessary Branches
- ③ Decrease Computation
- ④ less Time
- ⑤ Avoid Overfitting

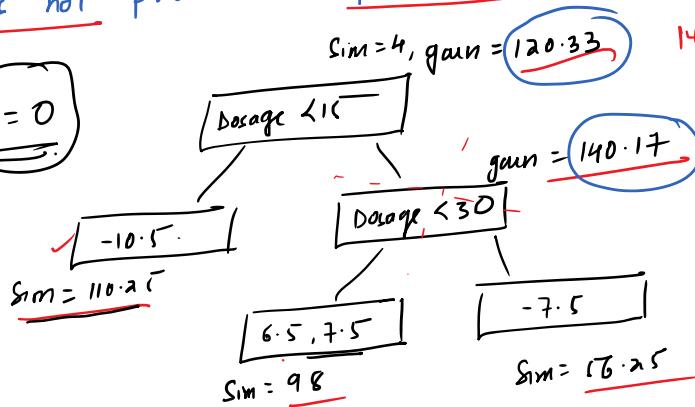
Pruning of Tree in xgBoost is purely done on gain.

Let us Assume a Threshold gain = γ (Gamma).

[Rule for Pruning if at Branch gain - $\gamma < 0$ then Prune and go above
if gain - $\gamma > 0$ then do not prune
+ if child is not pruned then parent also will not be pruned]

$$\text{if } \lambda = 0$$

$$\gamma = 130$$



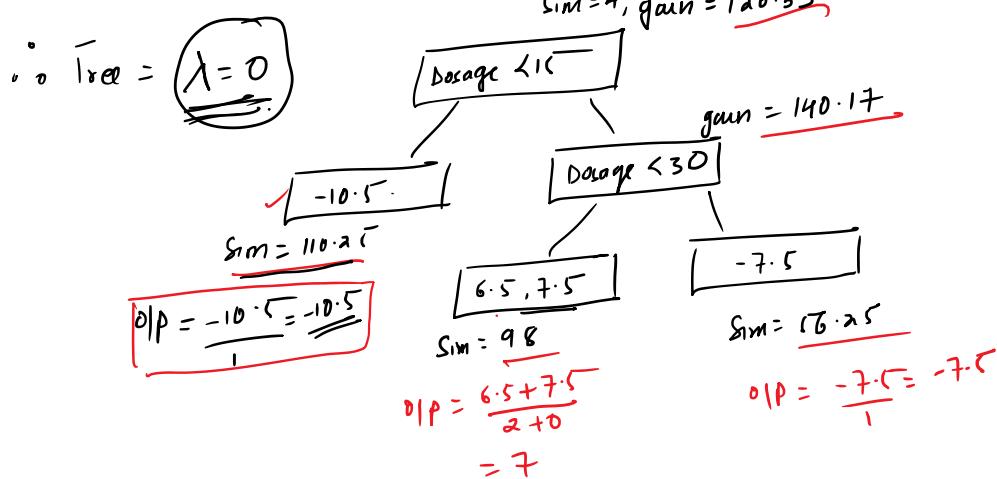
Here $120.33 - 130 < 0$ case for pruning But since child is not pruned so even root is not.
Here $140.17 - 130 = 10.17 > 0$

No pruning

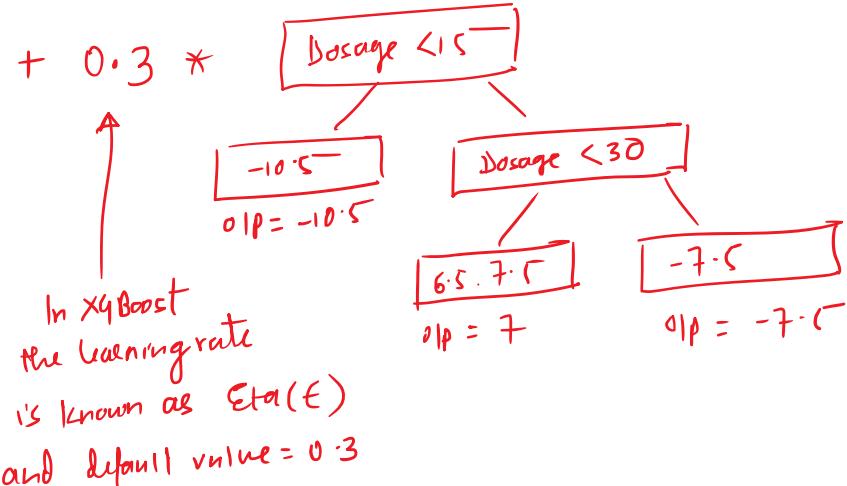
Here if γ is sufficiently large then it may result into pruning of tree.

Complete Tree till root \rightarrow known as Extreme Pruning.]

\rightarrow Regardless of value of λ & γ let us assume the tree as \rightarrow
Find output value at each leaf = $\frac{\text{Residual}}{\text{No of Residual} + \lambda}$



$$\text{New Prediction} = 0.5 + 0.3 * \text{Leaf Value}$$



$$\text{New Prediction for Dosage} = 0.5 + 0.3 * (-10.5)$$

$$\Rightarrow 0.5 + 0.3 * (-10.5) \Rightarrow \underline{\underline{-2.65}}$$

$$\begin{aligned} \text{Original Prediction} &= \underline{\underline{0.5}} \\ \text{New Prediction} &= \underline{\underline{-2.65}} \end{aligned}$$

* Find New prediction for all the other 3 Dosages

- * This will give New Residual.
- * Using the New Residual Build a New Tree and make New Prediction that will give smaller residual.
- * Keep Building the Tree until the max no of Tree constructed or residual is significantly small.

but for above case there are 4 Trees $\underline{T_1, T_2, T_3, T_4}$

$$\boxed{\text{Final Prediction} = 0.5 + 0.3 \times T_1 + 0.3 \times T_2 + 0.3 \times T_3 + 0.3 \times T_4.}$$