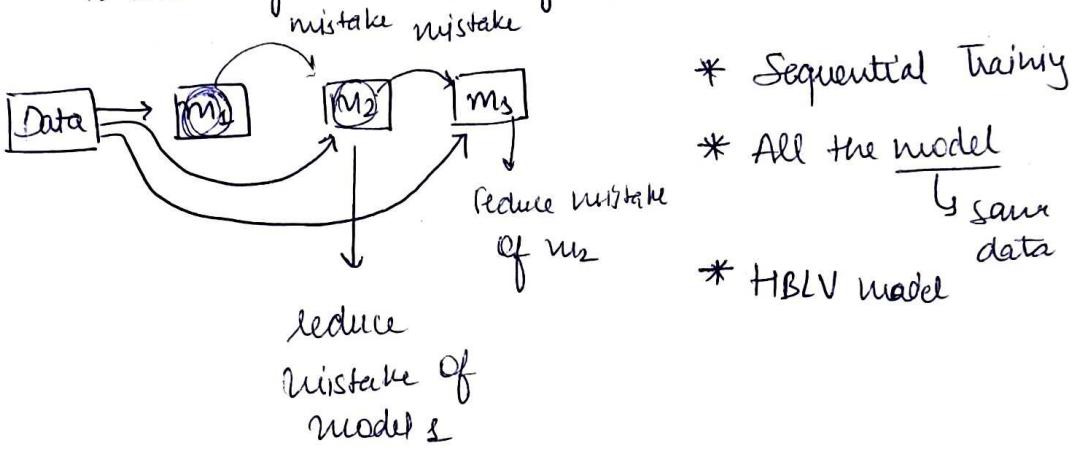


# Gradient Boosting

Boosting: Boosting is a general ensemble method in Machine learning that aims to create a strong classifier or regressor by combining the predictions of several weaker models. The idea is to build the strong model incrementally, by sequentially adding weak models that are trained to correct the mistakes made by the existing ensemble.



## What is Gradient Boosting

Gradient Boosting is a machine learning ensemble technique that aims to build a strong predictive model by combining the prediction of several weaker models using the concept of Additive Modelling typically decision trees. The method works by iteratively adding models to the ensemble, with each new model trained to correct the mistakes made by the combined ensemble of existing models.

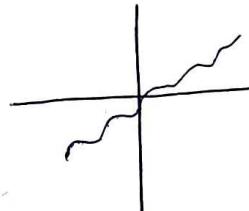
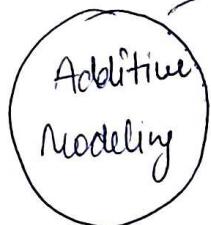
Gradient Boosting is a powerful and flexible method that can be used for both regression and classification tasks.

It is particularly effective when the data has complex, non-linear relationships, and it often performs well even with little hyperparameter tuning.

Its popularity in real-world application and machine learning competitions is testament to its effectiveness and its implementations in most major machine learning libraries, such as `sk-learn`, `XGBoost`, `LightGBM`,

`CatBoost`.

Gradient boosting work on mathematical principle



$$y = x + \sin x$$

$$y = x + y_{\text{bias}}$$

$$\boxed{w_1} \rightarrow \boxed{w_2} \rightarrow \boxed{w_n}$$
$$f_1(x) \quad f_2(x) \quad f_n(x)$$

$$f_1(x) + f_2(x) + f_3(x)$$

How Gradient Boosting Works? (Regression)

X		Y
---	--	---

$$y = f(\text{in})$$

find small funt<sup>n</sup>

find the relationship beth x, y  
with the help of gradient  
boosting.

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

Steps: Initialize  $f_0(x) = \arg \min_f \sum_{i=1}^n L(y_i, f)$

Step 2: for  $m = 1$  to  $M$ :

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

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

(b) fit a regression tree to the target  $r_{im}$  giving terminal regions.

$$R_{jm}, j = 1, 2, \dots, J_m$$

(c) for  $j = 1, 2, 3, \dots, J_m$  compute

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

(d) Update  $f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} v_{jm} \mathbb{1}_{x \in R_{jm}}$

Steps: Output  $\hat{f}(x) = f_M(x)$ .

Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, f(x))$ , number of iterations  $M$ .

$$\hat{y} = f(x)$$

$L(y, \hat{y})$   
output prediction

$\{(x_i, y_i)\}_{i=1}^n \rightarrow \text{rows}$

Squared loss  $L = \sum_{i=1}^n (y_i - \hat{y}_i)^2$

regression  $L(y, \hat{y}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2$

$$\boxed{L(y, \hat{y}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

## Explain Step 1

$$f_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$$

$$\frac{\partial L}{\partial \gamma} = 0$$

no. of rows

$$\begin{aligned}\frac{\partial f_0(\gamma)}{\partial \gamma} &= \frac{\partial}{\partial \gamma} \sum_{i=1}^n (y_i - \gamma)^2 \\ &= -2 \sum_{i=1}^n (y_i - \gamma) = 0\end{aligned}$$

$$y_1 - \gamma + y_2 - \gamma + y_3 - \gamma = 0$$

find that type  
of value of gamma which decreases  
value of these

	RED	OPS	Marketing	Profit
0	165	137	472	192
1	101	92	250	144
2	29	127	201	91

$$192 - \gamma + 144 - \gamma + 91 - \gamma = 0$$

$$3\gamma = 192 + 144 + 91$$

$$\gamma = \frac{192 + 144 + 91}{3}$$

$\gamma$  mean

Explain Step 2: As many decision tree  $\Rightarrow$

No. of loops 2nd  
step will run.

$m = \text{decision tree}$

$m_1 \rightarrow \text{first decision tree } (m=1)$

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

$$\gamma_{i1} = \left[ \frac{\partial L(y_i | f(1))}{\partial f(1)} \right]_{f=f_0} = \left[ \frac{\partial L(y_i | f_0(1))}{\partial f_0(1)} \right]$$

$$L(y, \hat{y}) = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \gamma_{i1} = \frac{1}{2} \frac{\partial}{\partial f_0(1)} \sum_{i=1}^n (y_i - f_0(1))^2$$

$$= -\frac{2}{2} (y_i - f_0(x_i)) \Rightarrow u_{i1} = f_0(x_1) - y_1$$

$$\delta_{11} = f_0(x_1) - y_1 = 142 - 192 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Value from table}$$

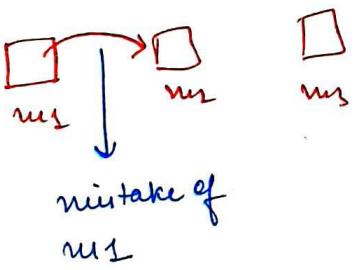
$$\delta_{21} = f_0(x_2) - y_2 = 142 - 144 \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

$$\delta_{31} = f_0(x_3) - y_3 = 144 - 91$$

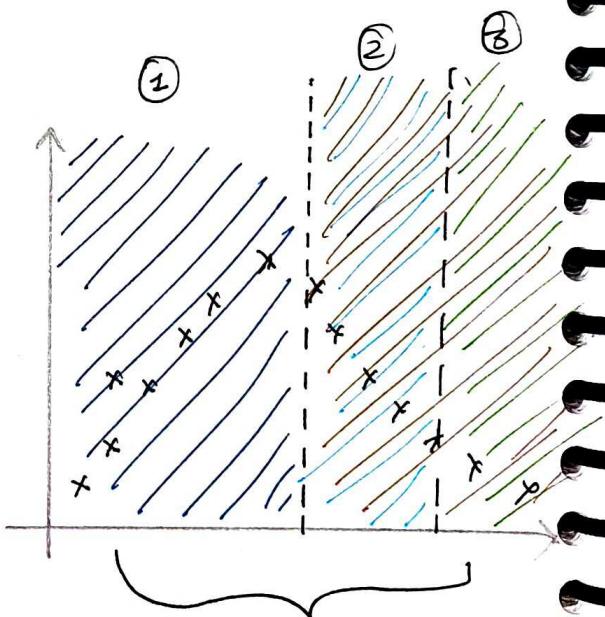
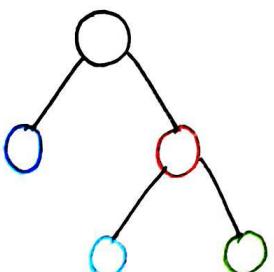
	R&D	Ops	Marketing	Profit	$f_0(x)$	R <sub>1</sub>
0	165	137	472	192	142.3	49.6
1	101	92	250	144	142.3	1.6
2	29	127	201	91	142.3	-51.3
.						<u>Output</u>

input

train decision tree as a input and output  $\rightarrow$  Step 2(b)



terminal region



terminal  
region

$$\gamma_{jm} = \operatorname{argmin}_r \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + r)$$

$$\gamma_{jz} = \operatorname{argmin}_r \left[ \sum_{x_i \in R_{jz}} L(y_i, f_0(x_i) + r) \right]$$

we have to find value

of  $r$  which help to reduce the value of

- \* we accept those  $x$  value which fall on specific region
- \* ex:-  $\gamma_{jm} = \gamma_{jz}$  we accept those  $x$  value which fall on first region of terminal region.

Step 2 (c)  $\rightarrow$

$$\gamma_{jL} = \operatorname{argmin}_r L(Y, f_0(x_3) + r) \quad (\because \text{not my } \sum \text{ bcz we my only one row})$$

$$\gamma_{jL} = \operatorname{argmin}_r (Y_3 - f_0(x_3) - r)^2 \quad (\because L = (y_i - \hat{y}_i)^2 \\ L = Y_3 - f_0(x_3) + r)^2$$

$$= \frac{\partial}{\partial r} \neq \operatorname{argmin}_r (Y_3 - f_0(x_3) - r)^2$$

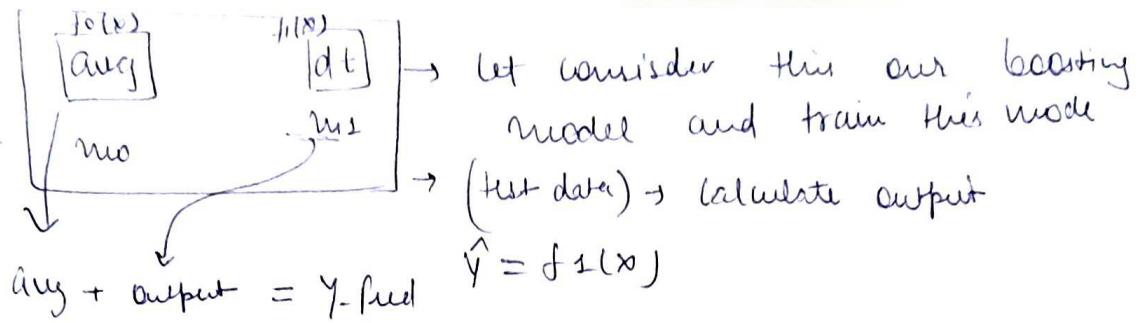
$$= -2(Y_3 - f_0(x_3) - r) = 0$$

$$Y_3 - f_0(x_3) - r = 0 \Rightarrow r = Y_3 - f_0(x_3) \\ = 91 - 142 \\ \boxed{r = -51}$$

Step 2 (d)  $\rightarrow$

$$f_m(x) = f_{m-1}(x) + \left[ \sum_{j=1}^{jm} \gamma_{jm} I(x \in R_{jm}) \right] \quad \nearrow m^2 \text{ output}$$

$$f_1(x) = f_0(x) + \text{output of first decision tree}$$



$$f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J^m} \gamma_{jm} I(x \in R_{jm})$$

$$f_5(x) = \underline{f_4(x)} + dt^5$$



$$f_4(x) = \underline{f_3(x)} + dt^4$$



$$f_3(x) = \underline{f_2(x)} + dt^3$$



$$\underline{f_1(x) + dt^2}$$



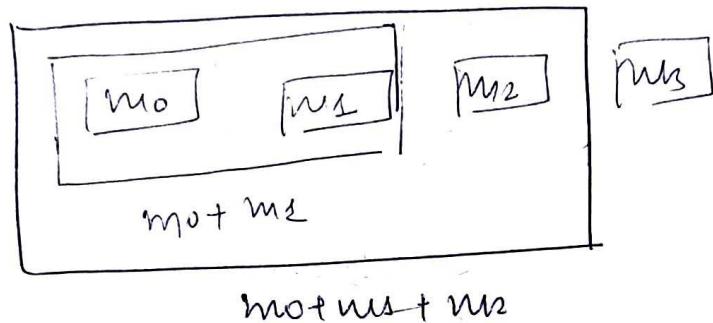
$$f_0(x) + dt^1$$

$$f_5(x) = f_0(x) + dt_1 + dt_2 + dt_3 = - dt_5$$

The what?

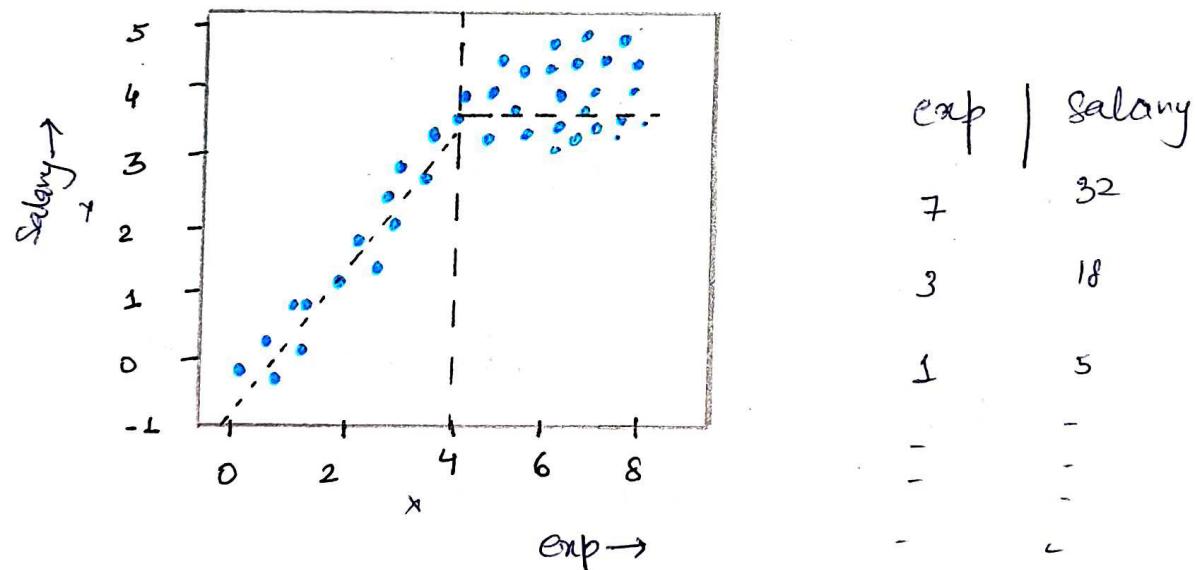
	Ren	Ops	Marketing	Benefit	$f_0(x)$	$y_1$	$f_1(x)$	$y - f_1(x)$
0	165	137	472	192	192	-	-	-
1	101	92	250	144	144	-	-	-
2	29	127	203	92	142	-	-	-

mean ↓ ↑



$$m_0 + m_1 + m_2 + m_3$$

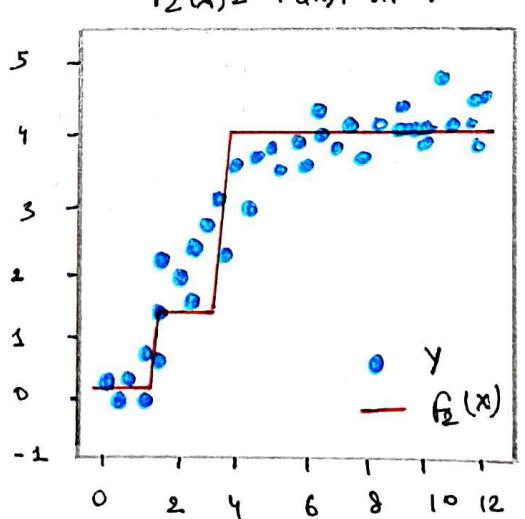
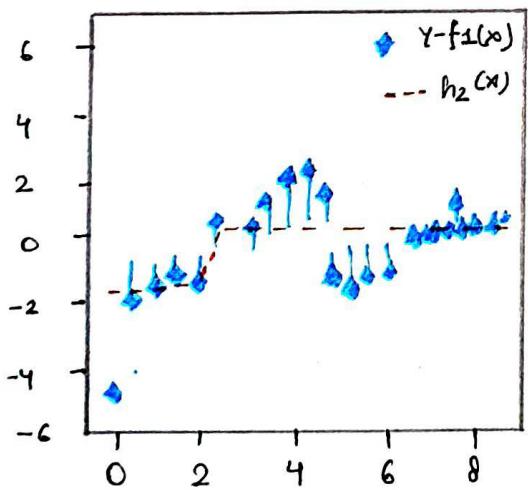
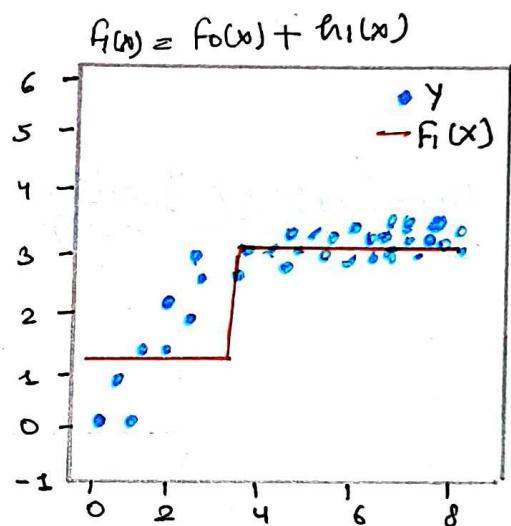
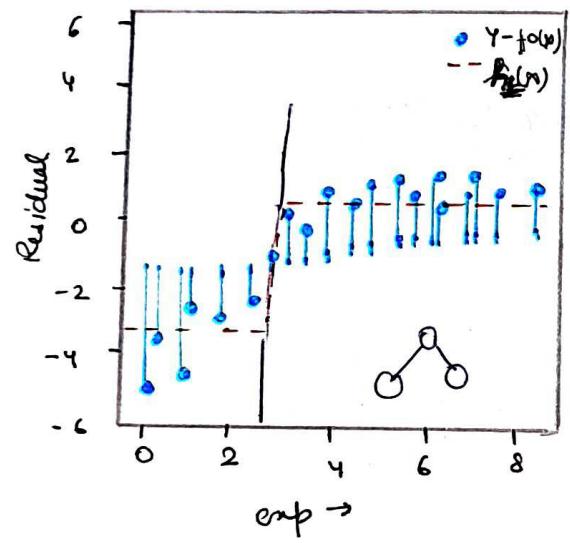
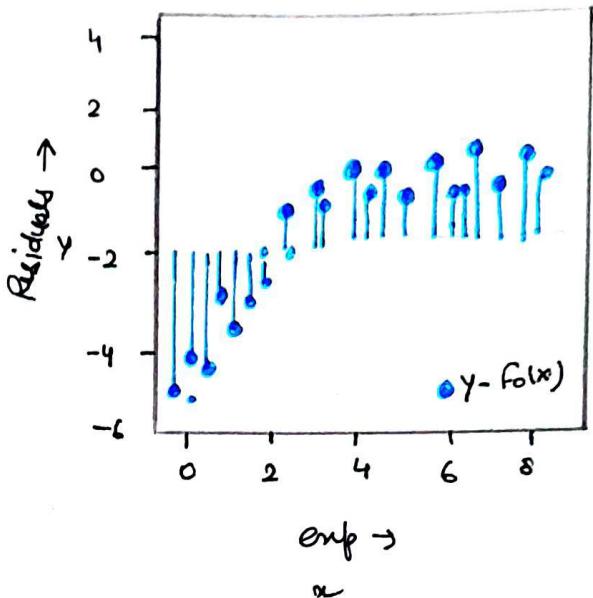
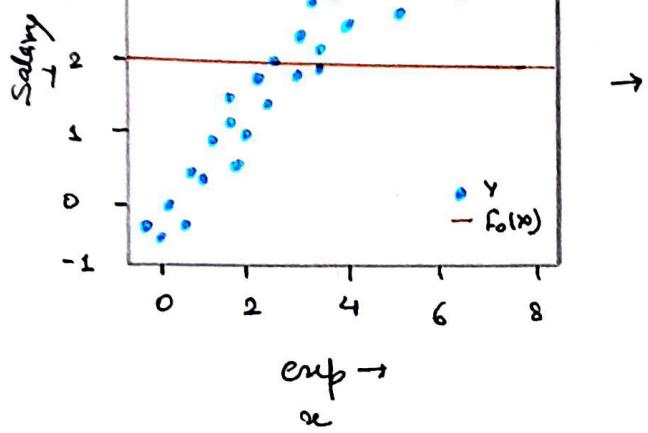
How Gradient Boosting Works for Regression?



we have to find

$$y = f(x)$$

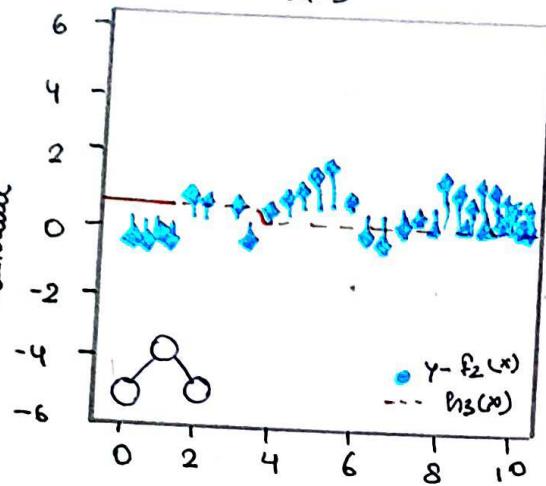
mean salary



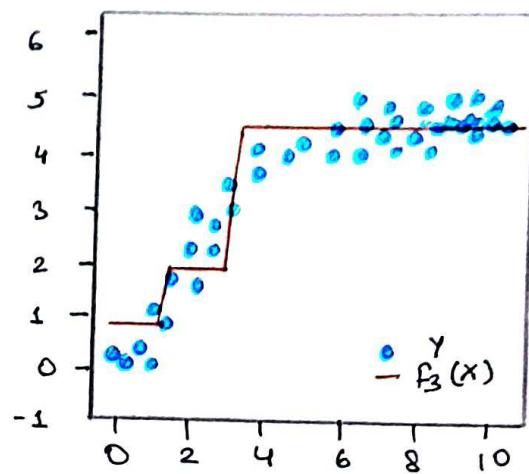
$$f_2(x) = f_0(x) + h_1(x) + h_2(x)$$

$m=3$

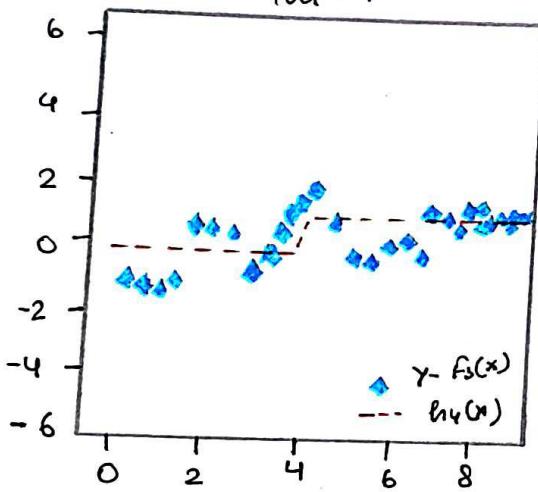
residual



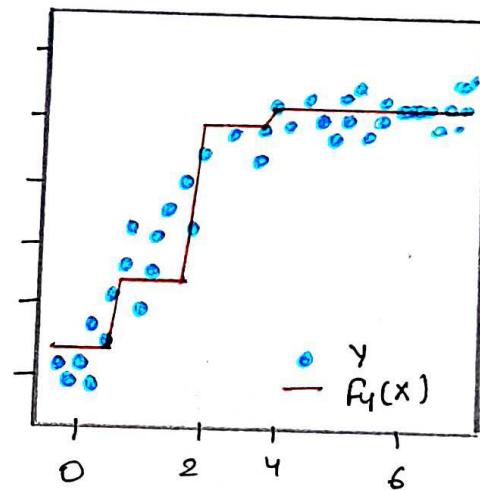
$$f_3(x) = f_0(x) + b_{1L}(x) + b_{1R}(x) + b_3(x)$$



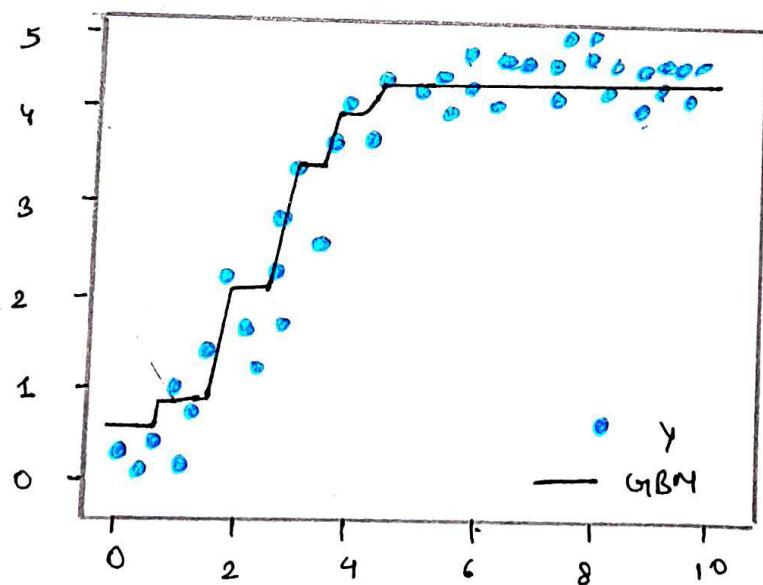
$m=4$



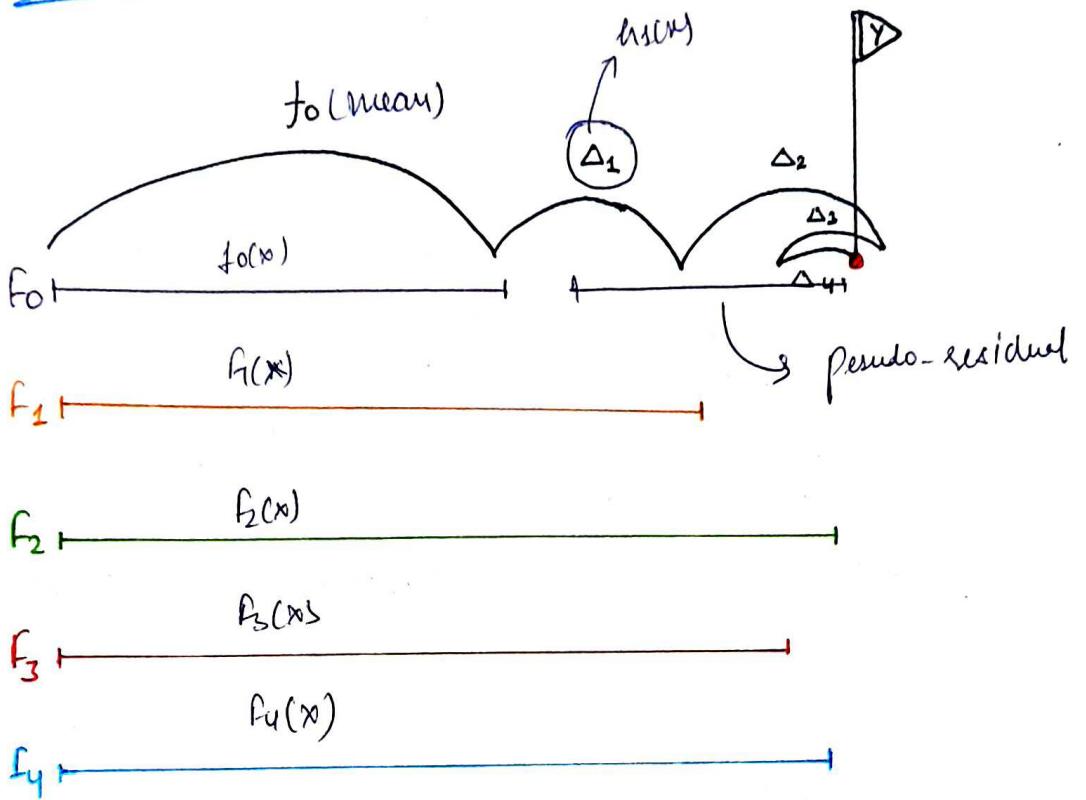
$$f_4(x) = f_0(x) + b_{1L}(x) + \dots + b_4(x)$$



10 Decision Tree

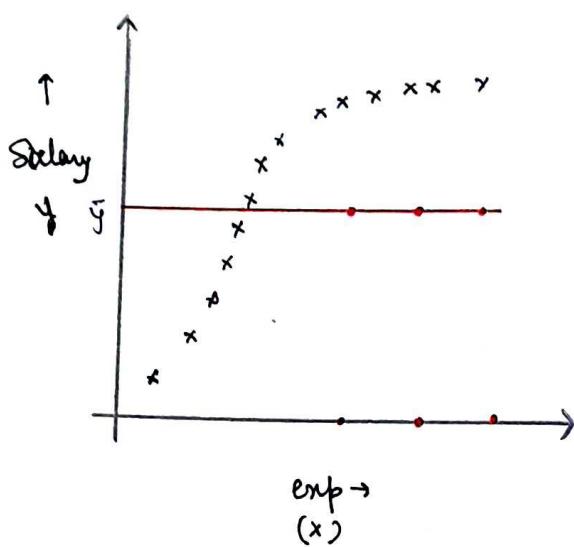


## Intuition



"Gradient Boosting is performing Gradient Descent in function space"

## Function Space Vs Parameter Space



emp | Salary | pred

10	(10)	(20)	} bcz mean output is all one same
20	(20)	(20)	
30	(30)	(20)	

MSE

$$L(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \frac{1}{3} [(10-20)^2 + (20-20)^2 + (30-20)^2]$$

## function Space

\* with  $\frac{1}{n}$  is mean square error

$$L(y, \hat{y}) = \left( \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \right) * \text{remove } \frac{1}{n} \text{ is squared error}$$

$$L(y, \hat{y}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \dots + (y_n - \hat{y}_n)^2$$

↑

constant  
(always present in data)

Variable  
(may vary)  
depend on model function

$$L(\hat{y}_1, \hat{y}_2, \hat{y}_3, \dots, \hat{y}_n)$$

math func

how many dims  
↓  
 $n+1$

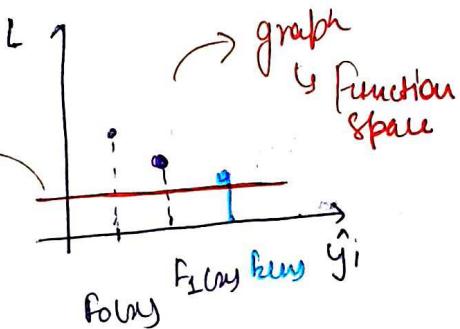
Loss function  $\hat{y}_i$  (n items)



$n+1$  dimension

relationship

\* we cannot draw  $n+1$  dimension graphs, so, we assume 2-D graph



we have to find  
the func where  
loss is low  
(True function)

$$f_{\text{true}}(x) = \bar{y} = 20$$

$$\hat{y}_i = f_{\text{loss}}(x) = \bar{y} = 20$$

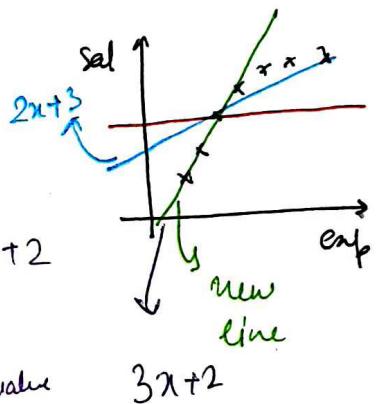
$$f_2(x_i) = \hat{y}_i = 2x_i + 3$$

↳ new loss value

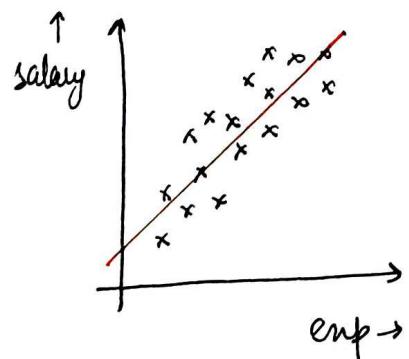
$$f_1(x_i) = \hat{y}_i = 3x_i + 2$$

↳ new value

↳ new loss value

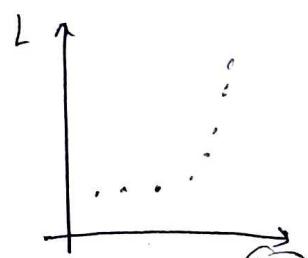


parametric  
alg → linear reg



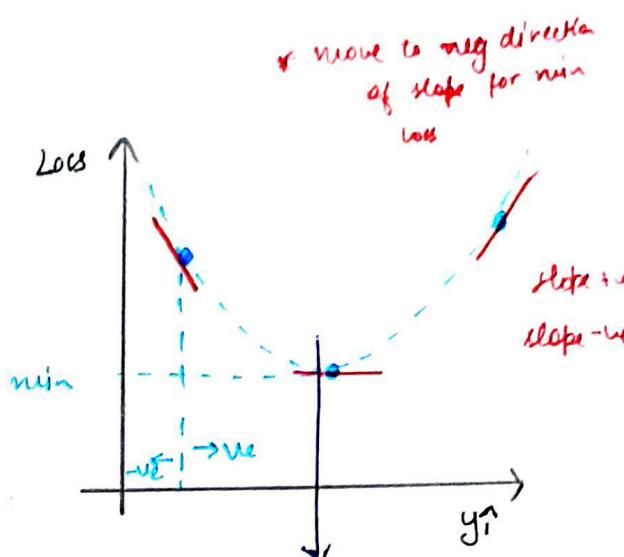
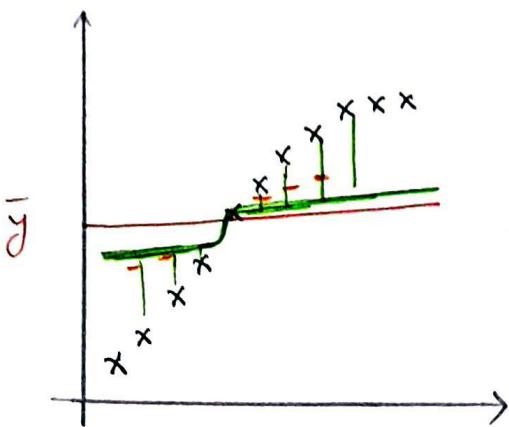
$$y = f_{\text{true}} = mx + b$$

↑  
m, b      logn  
              sine



\* In parametric algo, we already know which func is used. We have to just change the params and compare with loss function.

## 2. Direction of loss minimization



\* This whole step are doing in end step (a)  
(How gradient boosting works)

$$\underset{\text{I Decision Tree}}{\lim_{\downarrow}} = \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}} \quad \begin{matrix} \text{solve} \\ \boxed{y - f(x_i)} \rightarrow \text{step} \end{matrix}$$

This is ND curve in real so, we have to find n slope

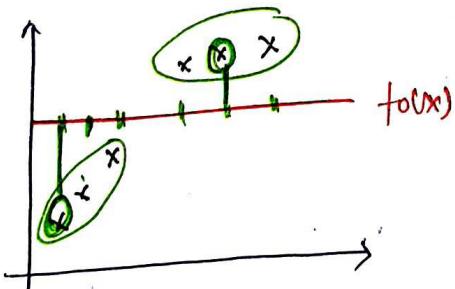
$$\left[ \frac{\partial L}{\partial y_1}, \frac{\partial L}{\partial y_2}, \dots, \frac{\partial L}{\partial y_n} \right]$$

help to find the direction.

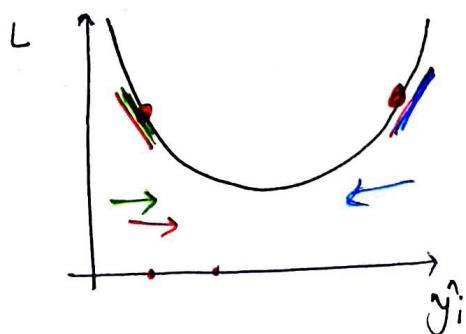
$$\rightarrow y_i - f(x_i)$$

$$\rightarrow y_2 - f(x_2)$$

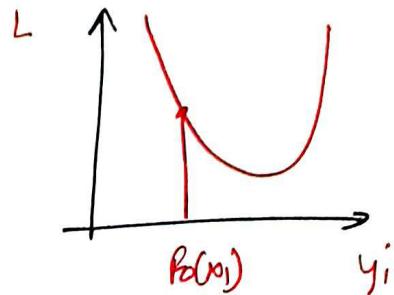
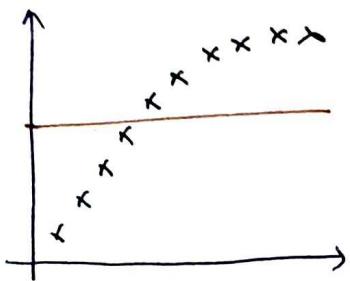
$$\rightarrow y_m - f(x_m)$$



$$\boxed{y_i - f(x_i)} \rightarrow \text{increase (+ve)} \\ \rightarrow \text{decrease (-ve)}$$



### 3. Update the function



$$-\left[ \frac{dL}{dy_i} \right] \rightarrow \text{direction}$$

$\downarrow$

$$\boxed{|y_i - f_0(x_i)|}$$

$$f_1(x_i) = f_0(x_i) - \text{slop}$$

$$f_1(x_i) = f_0(x_i) + y_i - f_0(x_i)$$

$$f_{\text{old}} = \text{m old}$$

$$f_1(x_i) = m_{\text{new}}$$

$$f_1(x_i) = f_0(x_i) + y_i - f_0(x_i)$$

$$\boxed{f_1(x_i) = y_i}$$

output of the data

( same output  
assign in  
new function )

~~massive overfitting~~

\* That's why we not use gradient descent method.

$$\text{slope} = y_i - f_0(x_i) = a$$

$$h_1(x_i) \approx a$$

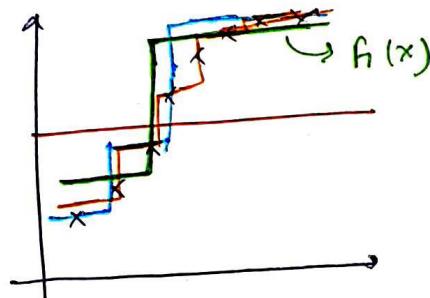
↳ approx equal to become if equal to  
a then again cancel the slope with old function.

$X$  | residual  $\rightarrow$  weak  
 $X$   $h_1(x)$   $\rightarrow$  predict weak or not perfect

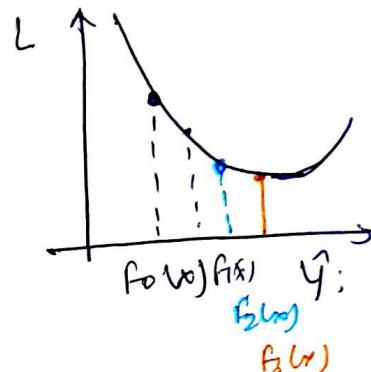
- \* If predict = not perfect then not cancel with old function.  
 ↳ not perfect but in right direction.

$$f_1(x) = f_0(x) + h_1(x)$$

#### 4. Iterate



$$f_1(x) = f_0(x) + h_1(x)$$



$$X \rightarrow \text{residual} \rightarrow h_2(x)$$

$$f_2(x) = f_0(x) + h_1(x) + h_2(x)$$

$$X \rightarrow \text{residual} \rightarrow h_3(x)$$

$$f_3(x) = f_0(x) + h_1(x) + h_2(x) + h_3(x)$$

## Difference Between Gradient Boosting and Gradient Descent

- 2) Gradient Descent → Parameter Space  
Gradient Boosting → Function Space

$$m_n = m_0 - \eta \text{ slope}$$

$$f_m(x) = f_{m-1}(x) + \eta h_m(x)$$

### Advantages

- \* Use any loss function, it give same result.  
mean squared error, mean absolute error, log loss.
- \* We can use any weak model instead of  $h_m(x)$   
decision tree.

## Gradient Boosting for Classification

We use same steps of Regression Gradient Boosting in classification. One difference is  $\rightarrow$  log loss instead of  $\text{MSE}$  for classification

cgpa	iq	is-placed
------	----	-----------

6.82	118	0
6.36	125	1
5.39	99	1
5.50	106	1
6.39	148	0
9.13	148	1
7.17	147	1
7.72	72	0

$$\begin{array}{ccc} \square & \square & \square \\ f_0(x) & f_1(x) & f_2(x) \\ \hline \end{array} \rightarrow \boxed{\text{Simple model}}$$

$$\log(\text{odds}) \rightarrow \log\left(\frac{\text{no. of } 1}{\text{no. of } 0}\right)$$

$$\begin{array}{c} \log\left(\frac{5}{3}\right) = f_0(x) \\ \log e \quad \downarrow 0.51 \end{array} \rightarrow \text{1st model}$$

cgpa	iq	is-placed	<u>pre1(log-odds)</u>
------	----	-----------	-----------------------

6.82	118	0	0.51
6.36	125	1	0.51
5.39	99	1	0.51
5.50	106	1	0.51
6.39	148	0	0.51
9.13	148	1	0.51
7.17	147	1	0.51
7.72	72	0	0.51

\* is-placed is kind of probability. So, we cannot combine is-placed with pre1(log-odds). We have to find probability of column pre1(log-odds).

We use the formula

$$P = \frac{1}{1 + e^{-\text{log-odds}}} = \frac{1}{1 + e^{-0.51}} = 0.62$$

Cgpa	iq	is-placed	pred1(log-odd)	pred (probability)
6.82	118	0	0.51	0.62 → 1
6.36	125	1	0.51	0.62 → 1
5.39	99	1	0.51	0.62 → 1
5.50	106	1	0.51	0.62 → 1
6.39	148	0	0.51	0.62 → 1
9.13	148	1	0.51	0.62 → 1
7.17	147	1	0.51	0.62 → 1
7.22	72	0	0.51	0.62 → 1

+0.62

→ In classification prob output is 0 or 1

let assume threshold is 0.5

$\frac{\uparrow 1}{\downarrow 0}$

\* calculate pseudo residual

$\text{res1} = \text{output} - \text{prediction} (\text{pred probability})$

X

Y

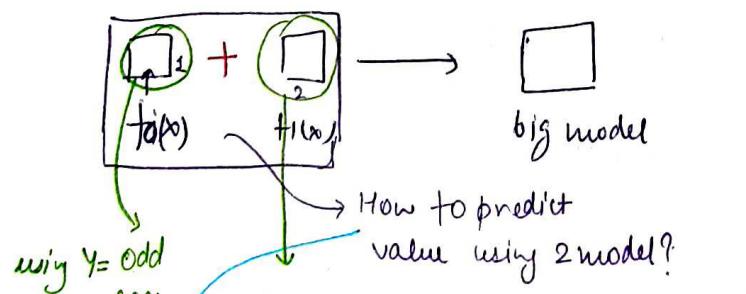
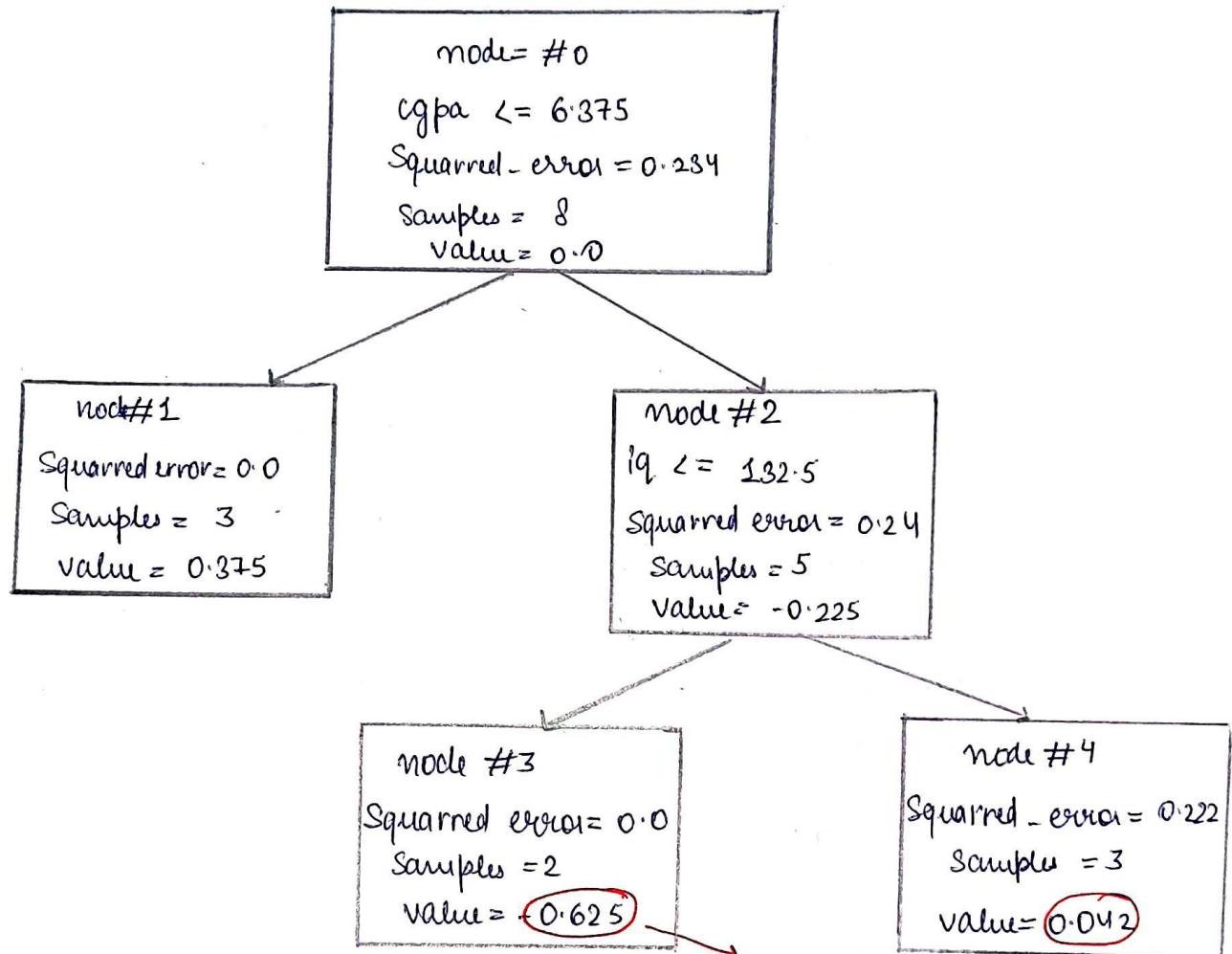
Cgpa	iq	is-placed	pred1(log-odd)	pred (probability)	res1
6.82	118	0	0.51	0.62	-0.62
6.36	125	1	0.51	0.62	0.37
5.39	99	1	0.51	0.62	0.37
5.50	106	1	0.51	0.62	0.37
6.39	148	0	0.51	0.62	-0.62
9.13	148	1	0.51	0.62	0.37
7.17	147	1	0.51	0.62	0.37
7.22	72	0	0.51	0.62	-0.62

→ Now, we can apply decision Tree on X and Y  
 where X is Cgpa and iq and Y is res1

→ Use decision tree on X and Y.

↳ regression tree bcz Y is number  
using regression tree in classification

weak learner → leaf node = 3



→ We have to add the prob of model 1 and model 2  
↳ diff of probability

In regression tree value are output but in classification these values are prob.  
bcz we train our model on differences of prob.

→ And we cannot add these prob with log odds. So, we have to convert prob into log odds.

$$\text{log(odd)} = \frac{\sum \text{Residuals}}{\sum (\text{Previous Prob} * (1 - \text{Previous Prob}))}$$

→ use this formula to every node to change value from ~~prob~~ probability into logodds.

cgpai	id	Is-placed	Res(log-odd)	Res(prob)	level	leaf-entropy
6.82	-	-	-	0.625	-0.625	3
	-	-	-	-	-	1
	-	-	-	-	-	1
	-	-	-	-	-	1
	-	-	-	-	-	4
	-	-	-	-	-	4
	-	-	-	-	-	4
	-	-	-	-	0.625 - 0.625	3

we find which row fallen on which leaf node using sklearn

- ex:- 1 row → 3 leaf node  
 2 row → 1 leaf node  
 5 row → 4 leaf node

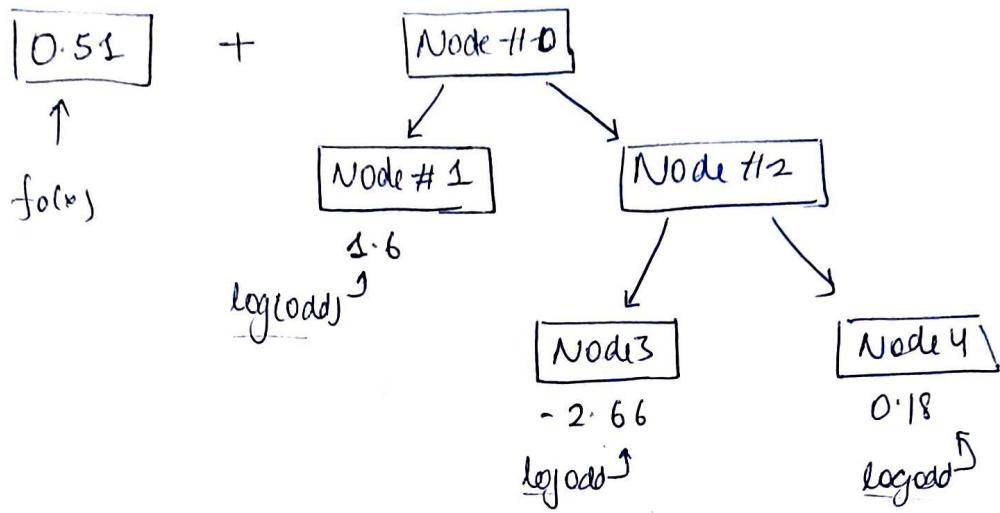
Let's calculate the log(odd) of node 3.

$$\text{log(odd)} = \frac{-0.625 - 0.625}{0.625(1 - 0.625) + 0.625(1 - 0.625)}$$

$$= \frac{-2 \times 0.625}{(0.625 \times 0.375) + (0.625 \times 0.375)}$$

$$= \frac{-2 \times 0.625}{2 \times 0.625 \times 0.375} = \frac{-1}{0.375} = -2.66$$

log odd of  
leaf node 3



Gpa	id	is-placed	pre1(log-odd)	pres(prob)	hess	leaf-entry 1	Pre2(log-odd)
5.82	118	0	0.51	0.625	-0.625	3	-2.15
-	-	-	-	-	-	-	2.11
-	-	-	-	-	-	-	2.11
-	-	-	-	-	-	-	2.11
-	-	-	-	-	-	-	0.69
-	-	-	-	-	-	-	0.69
-	-	-	-	-	-	-	0.69
7.72	72	0	0.51	0.625	-0.625	3	-2.15

$$\text{Row 1} = 0.51 + (-2.66) \xrightarrow{\text{3 leaf nodes}} -2.15$$

$\downarrow$        $\downarrow$   
 $f(x)$        $dt(\log \text{odd})$

\* Now again convert ~~Pre2(log-odd)~~ into prob to find out residual(s)

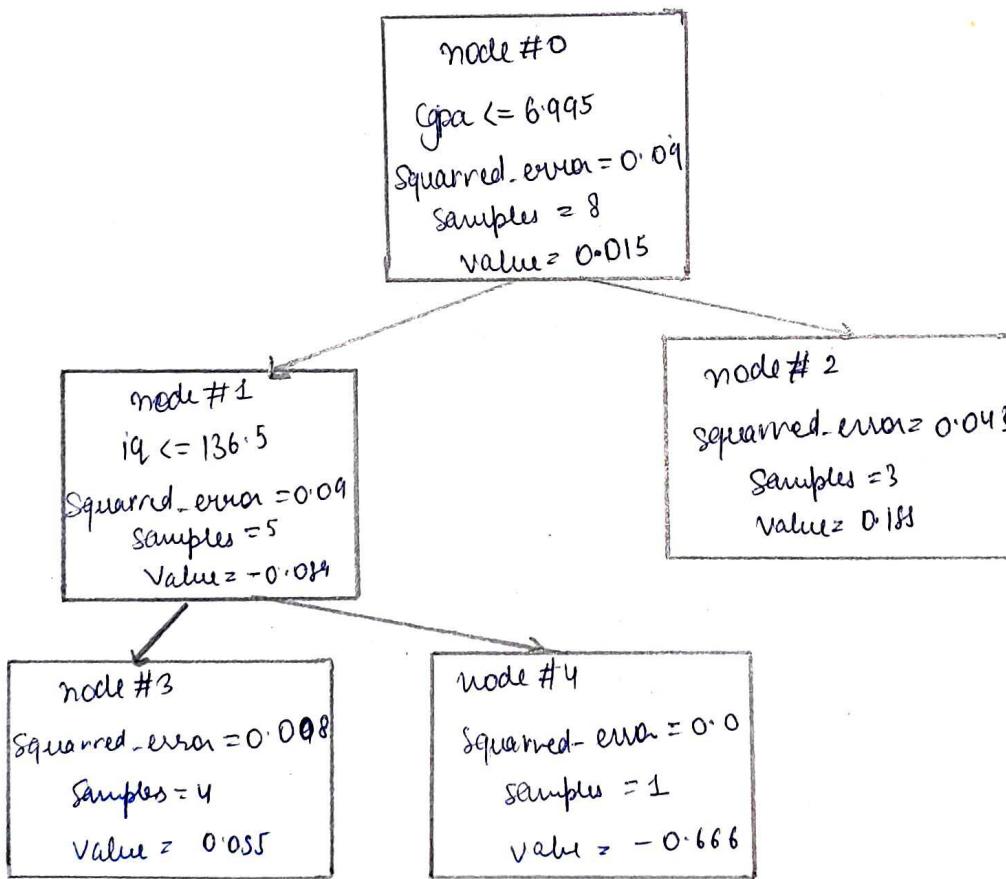
$$\text{probab (Pre2(log-odd)) Prob} = \frac{1}{1+e^{-\text{log odd}}}$$

cgpa	iq	is\_place	res1(log-odd)	res1(prob)	res1	Leaf\_entry1	pre2(log\_odd)	pre2(prob)	[res2]	$\rightarrow Y$
									0.103472	-0.103472
									0.891451	0.108049
									0.39151	0.103049
									0.89151	0.108049
									0.666151	-0.666151
									0.666151	0.333849
									0.666151	0.333849
									0.103472	-0.103472

$$res2 = Is\_place - pre2(prob)$$

↳ residual error)

Now again train <sup>second</sup> decision Tree with X (cgpa, iq) and Y (res2)



\* If big difference between res1 and res2 then we can multiply with learning rate ( $\eta$ ).

Ex:-  $res1 = -0.625$        $res2 = -0.103$

Ans now  $\uparrow$  while finding  $pre2$  (log odds) we can add  $\eta \times res1$

$$pre2(\text{log odds}) = 0.51 + N(1.6)$$

(  $\downarrow 10^{-1}$  )

→ Again we have to find the log odd of second decision tree using this formula

$$\frac{\sum \text{Residual}}{\sum (\text{Previous Prob} * (1 - \text{Previous Prob}))}$$

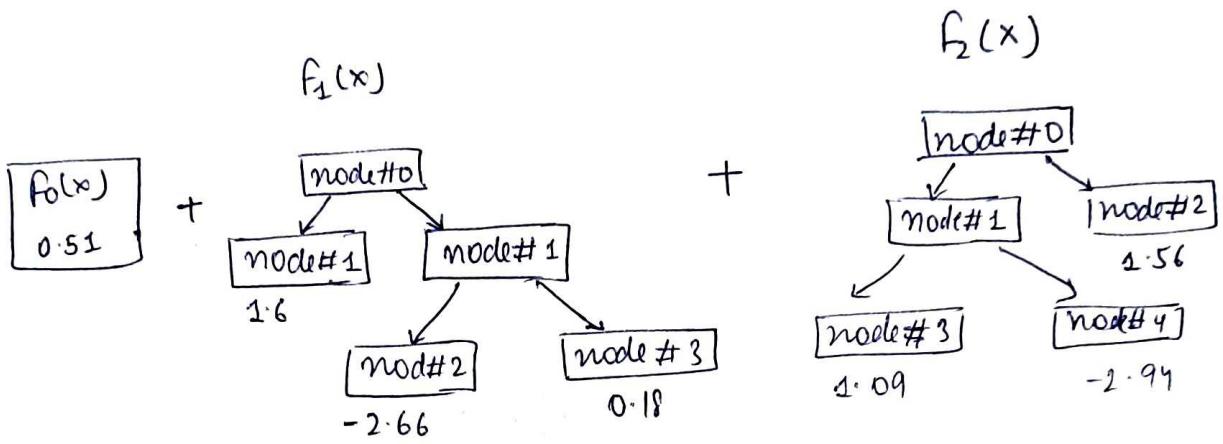
→ Again we have to find which row fallen from which node leaf.

Cgpa	Pre2(Prob)	Res2	leaf-entry
-	0.10	-0.1034	3
-	0.89	0.1080	3
-	0.89	0.1080	3
-	0.89	0.1080	3
-	0.66	-0.6661	4
-	0.66	0.3338	2
-	0.66	0.3338	2
-	0.103	-0.1034	2

Let's calculate the node 4 log odd value

$$\text{log odd} = \frac{-0.66}{0.66(1-0.66)} = -1 = -9.94$$

→ We can find all leaf node value using same step.



→ Add log odd of  $f_0(x)$  and  $f_1(x)$

$$f_{\text{pre}}(x) + \text{log odds of row (leaf entry)}$$

Ex:- Row 1  $\Rightarrow -2.15 + \frac{(3) \text{Node}}{1.09} = -2.15 + 1.09 = -1.06$

Row 2  $\Rightarrow 2.11 + \frac{3(\text{Node})}{1.09} = 2.11 + 1.09 = 3.20$

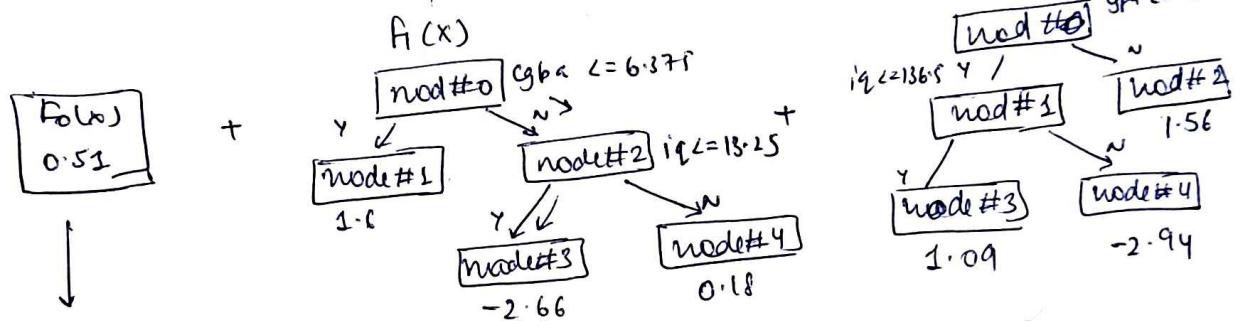
→ We have to find probability of log odd value.

so. we use  $P = \frac{1}{1+e^{-\text{log odds}}}$

lgrpa	leaf-entry <sup>2</sup>	pre3(log odds)	pre3(prob)
-	3	-1.068	0.155
-	3	3.201	0.960
-	3	3.201	0.960
-	3	3.201	0.960
-	4	-1.798	0.1420
-	2	2.516	0.904
-	2	2.316	0.904
-	2	-0.5983	0.354

## Prediction

new data = {7.2, 100y}  $\rightarrow$  {0, 1} y



$$0.51 + (-2.66) + 1.56 = -0.59 \rightarrow \text{log odds}$$

$$\text{prob} = \frac{1}{1+e^{-\text{log odds}}} = \frac{1}{1+e^{-0.59}} = 0.35$$

If threshold is 0.5 then output is 0.  
because prob < 0.5

## Maths Formulation

Step 0 → The loss function

$$\text{Log Loss} = L = -\frac{1}{n} \sum_{i=1}^n y_i \log p_i + (1-y_i) \log (1-p_i)$$

Gpa	iq	Placement	$\hat{y}_i   p_i$
8	80	1	0.62
7	70	0	0.32
6	60	1	0.51

$$L = -\frac{1}{n} \sum_{i=1}^n y_i \log p_i + (1-y_i) \log (1-p_i)$$

$$L = -\frac{1}{n} \sum_{i=1}^n y_i \log p_i + \log (1-p_i) - y_i \log (1-p_i)$$

$$= -\frac{1}{n} \left[ \sum_{i=1}^n y_i [\log p_i - \log (1-p_i)] + \log (1-p_i) \right]$$

$$= -\frac{1}{n} \left[ \sum_{i=1}^n y_i \log (\text{oddi}) + \log (1-p_i) \right]$$

↗ log odd  
 $\log \left( \frac{p_i}{1-p_i} \right)$

$$= -\frac{1}{n} \left[ \sum_{i=1}^n y_i \log (\text{oddi}) + \log \left( 1 - \frac{e^{\log (\text{oddi})}}{1 + e^{\log (\text{oddi})}} \right) \right]$$

$$p_i \geq \frac{e^{\log (\text{oddi})}}{1 + e^{\log (\text{oddi})}}$$

$$= -\frac{1}{n} \left[ \sum_{i=1}^n y_i \log (\text{oddi}) + \log \left( \frac{1}{1 + e^{\log (\text{oddi})}} \right) \right]$$

$$= -\frac{1}{n} \left[ \sum_{i=1}^n y_i \log(\text{odds}_i) + \log 1 - \log(1 + e^{\log(\text{odds}_i)}) \right]$$

$$\boxed{L = -\frac{1}{n} \left[ \sum_{i=1}^n y_i \log(\text{odds}_i) - \log(1 + e^{\log(\text{odds}_i)}) \right]}$$

↳ for entire data

for single data ↴

$$\boxed{L = -y \log(\text{odds}_i) + \log(1 + e^{\log(\text{odds}_i)})}$$

Step 1 → Finding  $f_0(x)$

→ loss function.

$$1. \text{ Initialize } f_0(x) = \underset{r}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, r)$$

$$L = -\frac{1}{n} \left[ \sum_{i=1}^n y_i \log(\text{odds}_i) - \log(1 + e^{\log(\text{odds}_i)}) \right]$$

$$\underbrace{L(y_i, \log(\text{odds}_i))}_{\hookrightarrow L(y_i, r)}$$

$$L = \underset{r}{\operatorname{argmin}} -\frac{1}{n} \left[ \sum_{i=1}^n y_i r - \log(1 + e^r) \right]$$

$$\frac{\partial L}{\partial r} = -\frac{1}{n} \left[ \sum_{i=1}^n y_i - \frac{e^r}{1 + e^r} \right] = 0$$

$$= -\frac{1}{n} \sum_{i=1}^n y_i + \frac{1}{n} \sum_{i=1}^n \frac{e^r}{1 + e^r} = 0$$

$$= -\text{avg} + \frac{n}{n} \frac{e^r}{1 + e^r} = 0$$

$$\frac{e^r}{1 + e^r} = \text{avg}$$

$$e^r = \text{Pavg} + \text{Per} \Rightarrow e^r - \text{Pavg}e^r = \text{Pavg}$$

$$e^r(1-\text{Pavg}) = \text{Pavg}$$

$$e^r = \frac{\text{Pavg}}{1-\text{Pavg}}$$

$$\boxed{r = \log \left( \frac{\text{Pavg}}{1-\text{Pavg}} \right)} \rightarrow f_0(x) = r = \log \left( \frac{\text{Pavg}}{1-\text{Pavg}} \right)$$

Cgpa, iq, is\\_placed

6.82	118	0
6.36	125	1
5.39	99	1
5.50	106	1
6.39	148	0
8.9	103	1
9.13	148	1
7.17	142	1
7.72	72	0

$$\text{Pavg} = \frac{5}{8}$$

$$\log \left( \frac{5/8}{1-5/8} \right) = \log \left( \frac{3/8}{5/8} \right)$$

In Gradient Boosting  
Classification

$$\text{we used } \frac{\log \left( \frac{y_i}{f(x)} \right)}{f_0(x)}$$

Step 2  $\rightarrow$  Pseudo Residuals

2. For  $m=1$  to  $M$ :

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

$$\delta_{im} = \left[ \frac{\partial L(y_i, f_m(x))}{\partial f_m(x)} \right]_{f_m = f_{m-1}}$$

$$L = -y \log(\text{odds}_i) + \log(1 + e^{\log(\text{odds})})$$

$f(x) \rightarrow f(x) + \log \text{odd}$

$$\frac{\partial L}{\partial \log(\text{odds})} = -y + \left[ \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}} \right] \rightarrow \text{prob output}$$

$$= -y + p = y - p = \frac{|y_i - p_i|}{\log \text{odd}}$$

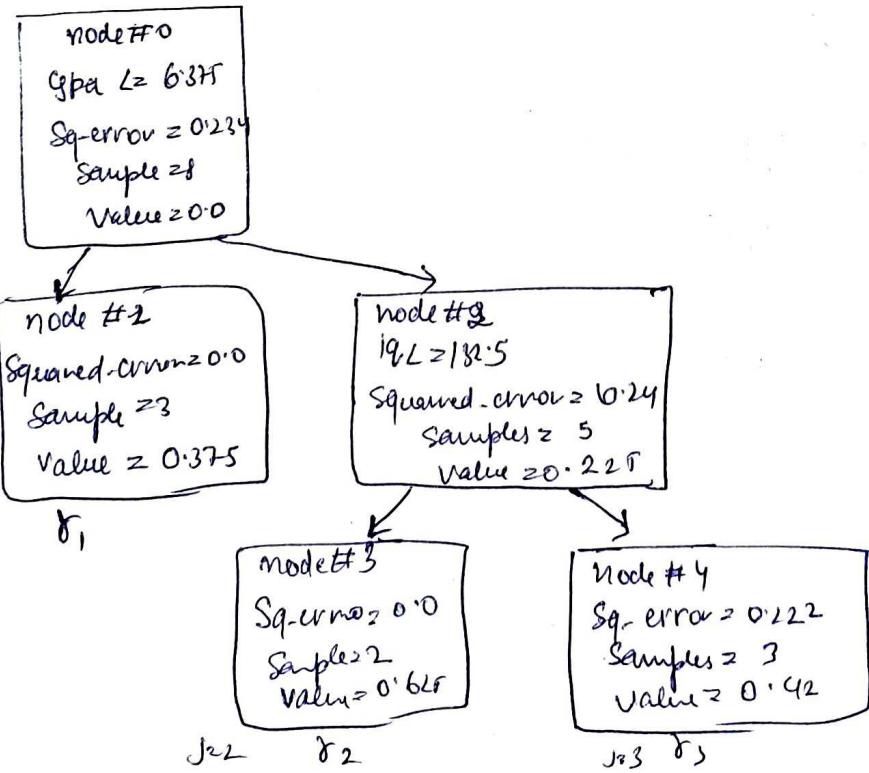
### Step 2.b → Train Regression Tree

There's no mathematical formula we have to just train on Regression Tree.

### Step 2.c → compute lambda for all leaf nodes

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

$$\gamma_{jm} = \operatorname{argmin}_{x_i \in R_{jm}} \sum L(y_i, f_{m-1}(x_i) + \gamma)$$



<u>(y<sub>fa</sub>)</u>	<u>iq</u>	<u>is placed</u>	<u>push(log-odds)</u>	<u>push(prob)</u>	<u>less</u>	<u>leaf-entrys</u>
					-0.62	3
					0.375	1
					0.375	1
					0.375	1
					-0.625	4
					0.325	4
					0.325	4
					-0.625	3

→ Find  $r$  for node leaf 2.

$$r_2 = \underset{r}{\operatorname{arg\,min}} \left[ L(y_1, f_0(x_1) + r) + L(y_2, f_0(x_2) + r) \right]$$

$$\frac{\partial L}{\partial r} = 0 \rightarrow \text{minimum value}$$

$$x \rightarrow f_0(x_1) + r$$

$$a \rightarrow f_0(x_1)$$

Taylor Series

$$f(x) = f(a) + f'(a)(x-a) + \frac{f''(a)(x-a)^2}{2!}$$

$$\Rightarrow L(y_1, f_0(x_1)) + \frac{\partial L}{\partial f_0(x_1)} r + \frac{1}{2} \frac{\partial^2 L}{\partial f_0(x_1)^2} r^2$$

$$\Rightarrow \frac{\partial L}{\partial y_1} = \frac{\partial L}{\partial f_0(x_1)} + \frac{1}{2} r \frac{\partial^2 L}{\partial f_0(x_1)^2} = 0$$

$$\Rightarrow \frac{r \frac{\partial^2 L}{\partial f_0(x_1)^2}}{\partial f_0(x_1)} = - \frac{\partial L}{\partial f_0(x_1)}$$

$$\delta = \frac{\frac{-\partial L}{\partial f_0(x_1)}}{\frac{\partial^2 L}{\partial f_0(x_1)^2}}$$

Numerator  $\rightarrow \frac{\partial L}{\partial f_0(x_1)} = \frac{-y_1 + e^{\log(\text{odds}_1)}}{1 + e^{\log(\text{odds}_1)}}$

\* Same step for 2.

denominator

$$\frac{\partial}{\partial f_0(x_1)} \frac{\partial L}{\partial f_0(x_1)} = \frac{\partial}{\partial \log(\text{odds}_1)} y_1 - \frac{e^{\log(\text{odds}_1)}}{1 + e^{\log(\text{odds}_1)}}$$

$$= \frac{\partial}{\partial \log(\text{odds}_1)} [y_1 - e^{\log(\text{odds}_1)} (1 + e^{\log(\text{odds}_1)})^{-1}]$$

$x \downarrow y$   
 $x' y + x y'$

$$= - \left[ e^{\log(\text{odds}_1)} (1 + e^{\log(\text{odds}_1)})^{-1} - e^{\log(\text{odds}_1)} (1 + e^{\log(\text{odds}_1)})^{-2} \right]$$

$$= \frac{e^{\log(\text{odds}_1)}}{(1 + e^{\log(\text{odds}_1)})^2} - \frac{e^{\log(\text{odds}_1)}}{(1 + e^{\log(\text{odds}_1)})}$$

$$= \frac{e^{\log(\text{odds}_1)}}{1 + e^{\log(\text{odds}_1)}} \left[ \frac{1}{1 + e^{\log(\text{odds}_1)}} - 1 \right] \rightarrow \frac{1}{1 + e^{\log(\text{odds}_1)}}$$

$$= \frac{e^{\log(\text{odds}_1)}}{(1 + e^{\log(\text{odds}_1)})} \left[ \frac{1}{(1 + e^{\log(\text{odds}_1)})} \downarrow \right] \rightarrow \frac{1}{1 + e^{\log(\text{odds}_1)}} \rightarrow P$$

$$\gamma = \frac{\text{residual}_i}{P_i(1-P_i)}$$

$$\gamma_2 = \frac{\sum \text{residual}_i}{\sum \text{Prob}_i(1-\text{Prob}_i)} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Probaed}$$

↳ by node 2

Step 2 D - Update the model

$$f_1(x) = f_0(x) + \text{output from dt of some leaf node}$$

$$f_2(x) = f_0(x) + f_1(x) + \text{output from dt+}$$

Step 3 :- Final Model

$$f_M(x) = f_{M-1}(x) + \text{last decision output}$$

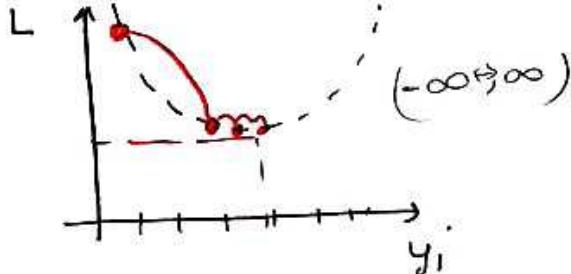
Boosting Model

## Log(Odds) vs Probability

Unconstrained Prediction Space: Log odds can span the entire real line ( $-\infty$  to  $\infty$ ), while probabilities are constrained between 0 and 1. Algorithms like gradient boosting involve adding correction via the weak learner to the prediction iteratively. If you're working in the log odds space, there's no need to worry about your predictions going out of bounds.

Better Gradients: When computing gradients (which guide the addition of new trees in boosting), the gradients can be more informative and have better magnitudes in the log odds space, than in the probability space, especially when probabilities are near 0 or 1.

log(Odds)  $\rightarrow$  Big jumps



prob  $\rightarrow$  small jumps

