

# XGBoost

## XGBoost Introduction

ML → [70's 80's] → Linear

Naive Bayes

90's → RF

SVM  
GB

Performance

↑  
overfitting

Scalability

↓  
speed

Not an

algorithm

XGBoost

ML + SF → library  
Gradient Boosting  
Software engineer

## Flexibility

1. Cross platform

linux

window

2. Multiple language support

Python

R

matlab

3. Integration with other libraries and tools

4. Support all kind of ML problems

Reg

Classification

Time series

Ranking

Binary

Multi

Python → API → JAVA

numpy / sklearn / pandas

model building

distributed → spark

model interpretable

model deployment

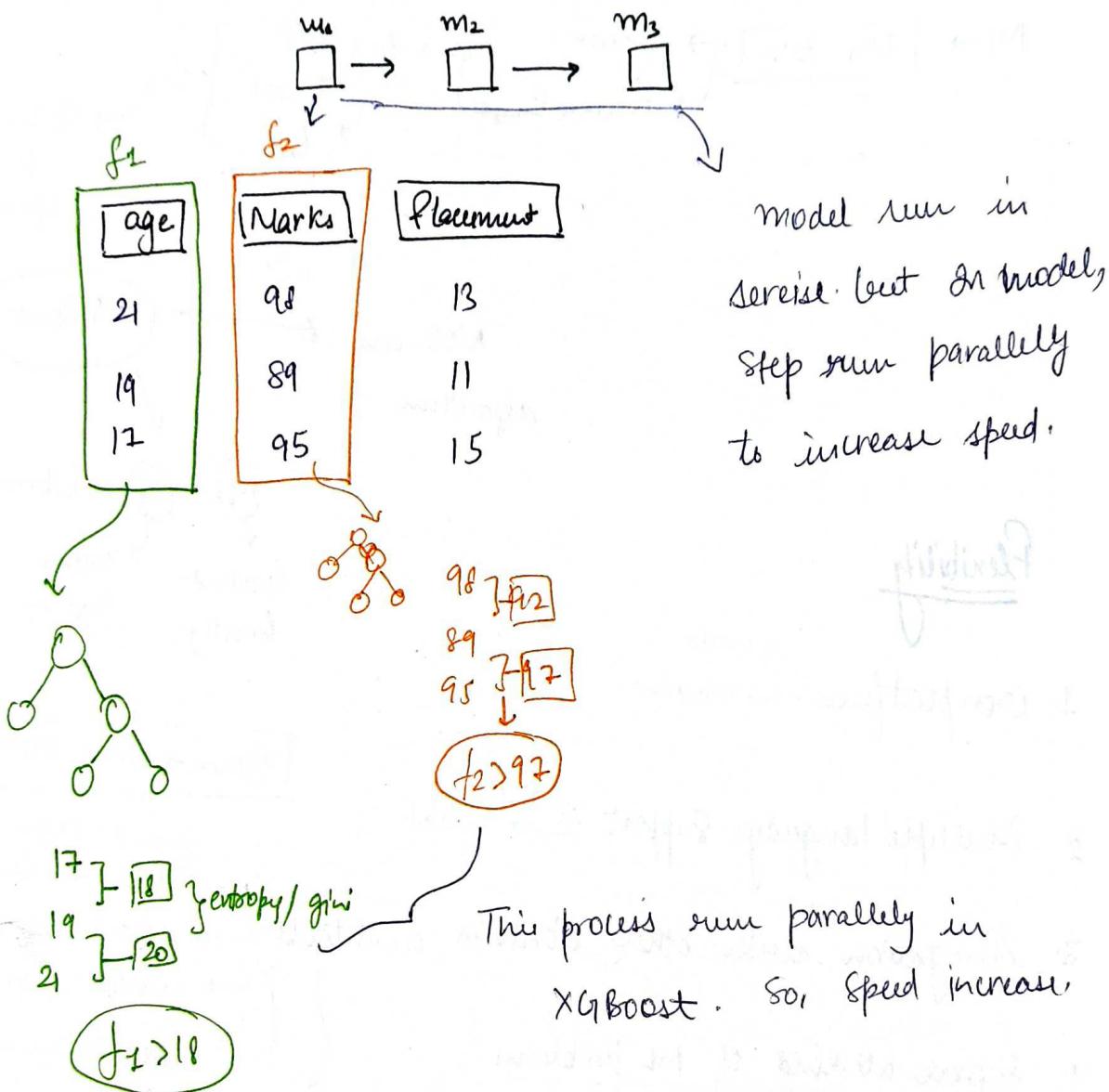
work flow

me flow

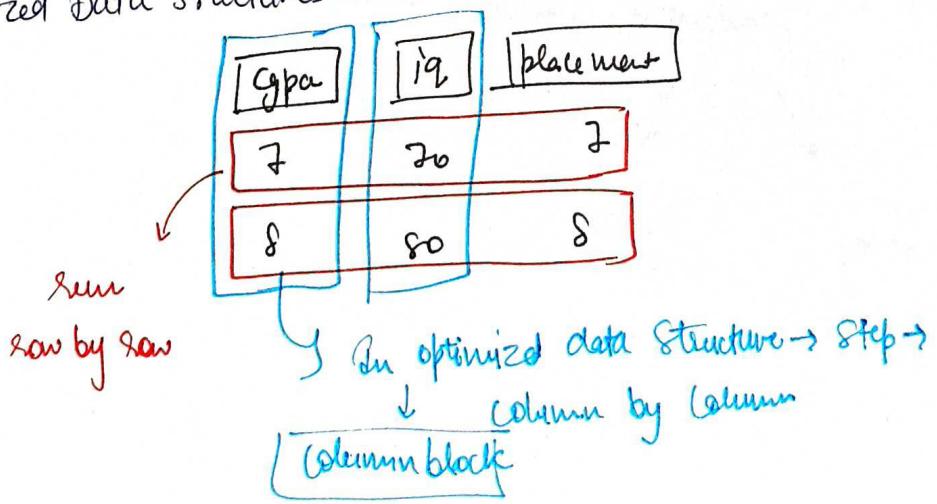
docker

## 2. Speed

### 1. Parallel Processing

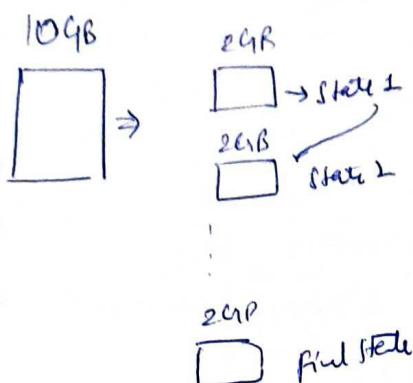


## 2. Optimized Data Structures

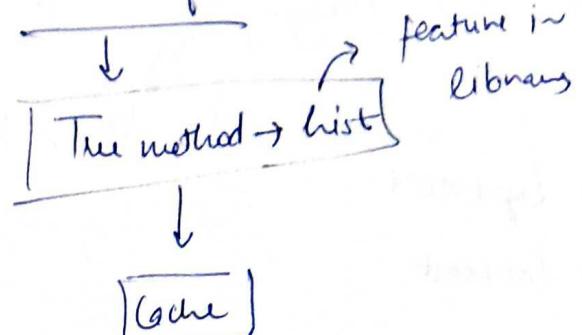


## 3. Cache Awareness

### 4. Out of Core Computing



XgBoost out of core



### 5. Distributed Computing

→ multiple machine

### 6. GPU Support

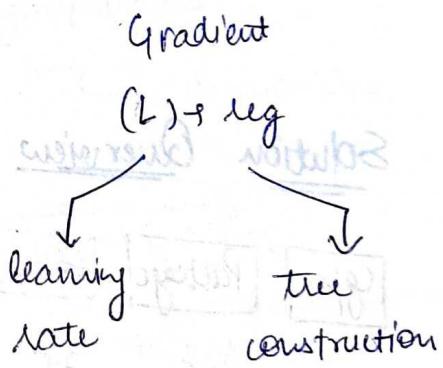
tree-method = gpu\_hist

## Performance

### 1. Regularized learning Objective

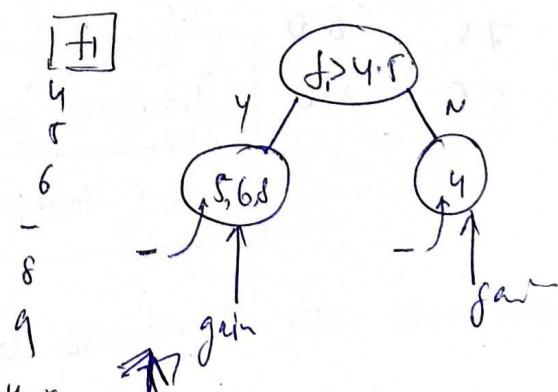
↳ linear reg  $\rightarrow (L_1 \& L_2)$

↳ Train  $\rightarrow$  Test  
✓ overfitting  $\times$



### 2. Handling Missing Values

### 3. Sparsity Aware split Finding



\* Missing value are the in Y and another time in N And then check gain

#### 4. Efficient Split finding (Weighted Quantile sketch + Approximate Tree learning)

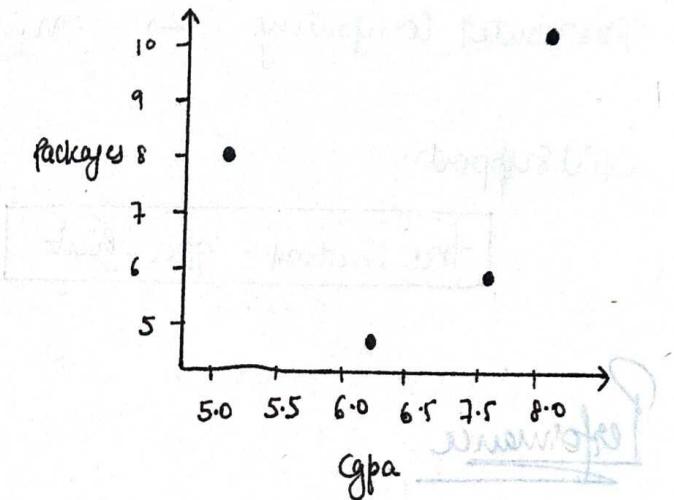
#### 5. Tree Pruning

↳ LightGBM

↳ CatBoost

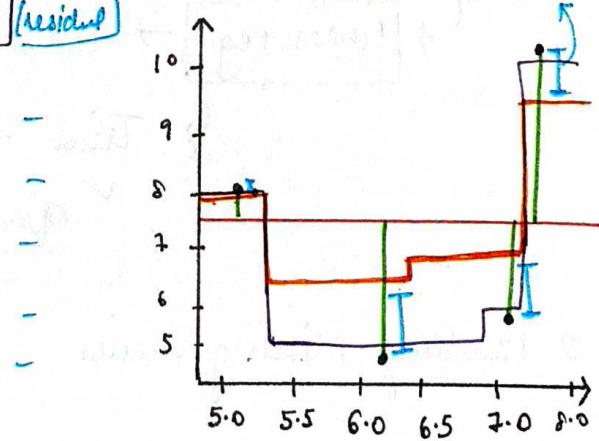
#### The Problem Statement

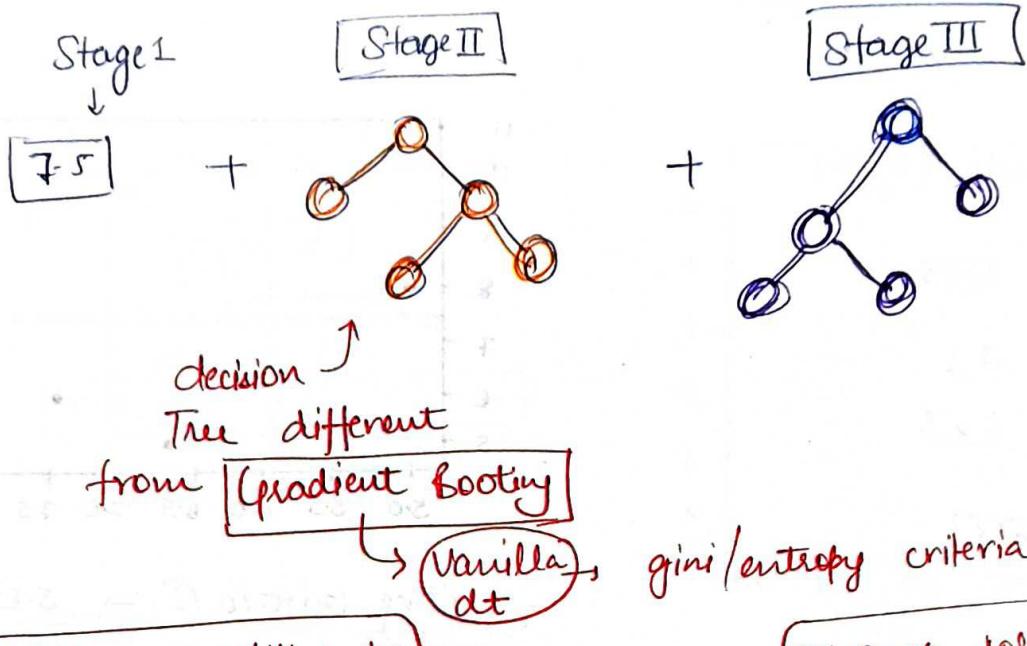
cgpa	Packages
6.7	4.5
9.0	11.0
7.5	6.0
5.0	8.0



#### Solution Overview

cgpa	Package	$f_0(x)$	residual	$f_1(x)$	residual
6.7	4.5	7.5	-3	-	-
9.0	11.0	7.5	3.5	-	-
7.5	6.0	7.5	-1.5	-	-
5.0	8.0	7.5	0.5	-	-





XGBoost use different criteria

gini/entropy criteria

→ XGBoost follow same flow as Gradient boosting but diff criteria in dt.

### Step by Step Calculation

Gpa	Package	mode1	residuals
6.7	4.5	7.3	-2.8
9.0	11.5	7.3	3.7
7.5	6.0	7.3	-1.3
5.0	8.0	7.3	0.7

Gpa | residuals

↓ decision tree

leaf node

↑ [-2.8, 3.7, -1.3, 0.7]

Calculate quantity to leaf node is

\* Same as we calculate gini impurity in Gradient boosting of DT.

← Similarity Score in XGBoost

$$SS = \frac{(\text{sum of residual})^2}{\text{no. of residual} + \lambda}$$

Let  $\lambda = 0$

$$= \frac{(-2.8 + 3.7 + 1.3 + 0.7)^2}{4 + 0}$$

= 10.2 → residual similar to each other score

\* We need to split residual with the help of cgpa column.

Ascending order (cgpa)

$$\begin{array}{l} 5.0 \\ 6.7 \\ 7.5 \\ 9.0 \end{array} \left[ \begin{array}{l} 5.85 \\ 7.1 \\ 8.25 \end{array} \right]$$

[cgpa] [residual]

5.0	0.7
6.7	-2.8
7.5	-1.3
9.0	3.7

Splitting criteria ①  $\rightarrow 5.85$

[cgpa < 5.85]

Splitting criteria ①  $\rightarrow 5.85$

[cgpa < 5.85]

0.7

-2.8, -1.3, 3.7

$$SS_L = \frac{(0.7)^2}{1+0} = 0.49$$

$$SS_R = \frac{(-2.8 + -1.3 + 3.7)^2}{3+0} = 0.05$$

$$gain = (SS_L + SS_R) - SS_{root}$$

$$= (0.49 + 0.05) - 0.02$$

= 0.52  
SS score increase

Splitting criteria ②  $\rightarrow 7.1$

[cgpa < 7.1]

$SS_L = \frac{(0.7 - 2.8)^2}{2} = 2.20$

$SS_R = \frac{(-1.3 + 3.7)^2}{2} = 2.88$

$$gain = (SS_L + SS_R) - SS_{root}$$

$$= 2.20 + 2.88 - 0.02$$

$$= 5.06$$

↑ SS score increases

Splitting criteria 3  $\rightarrow 8.25$

[cgpa < 8.25]

$SS_L = 3.85$

$SS_R = 13.69$

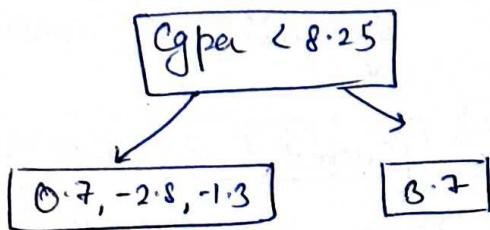
$$gain = (SS_L + SS_R) - SS_{root}$$

$$= 17.52$$

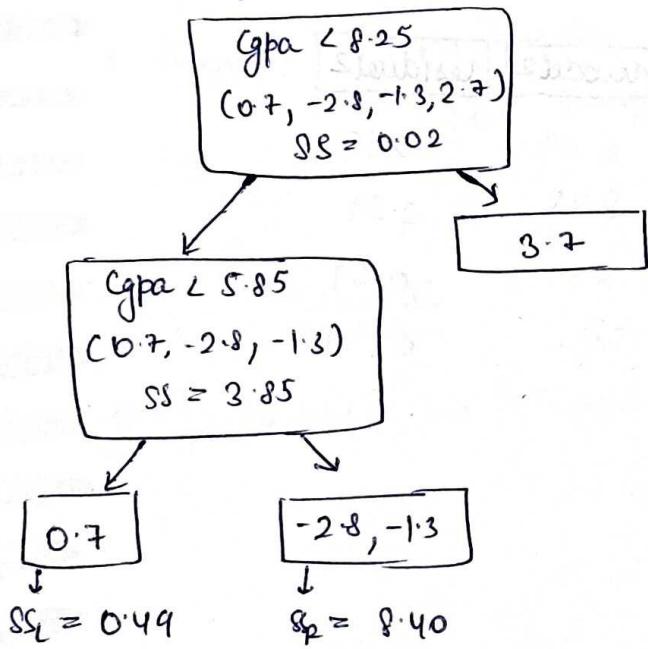
↑ SS score increase

→ Take Stage 2 because higher SS score increase than other  
Criteria 3

Stage 2 → decision tree



① Splitting Criteria  $\rightarrow 5.85$

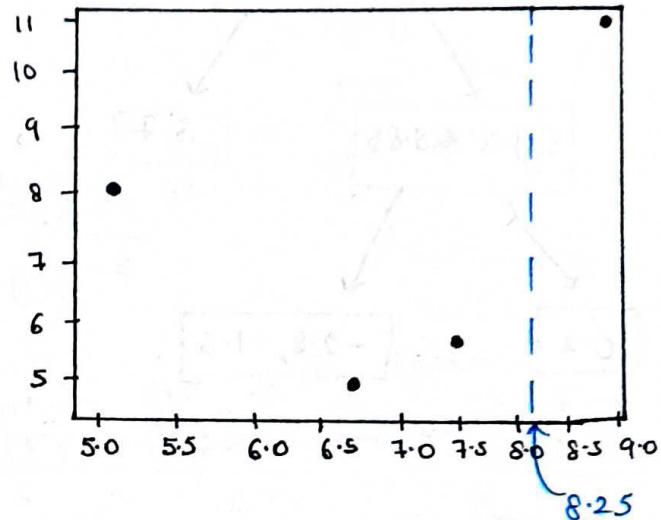


$$\text{gain} = (0.49 + 8.40) - 3.85$$

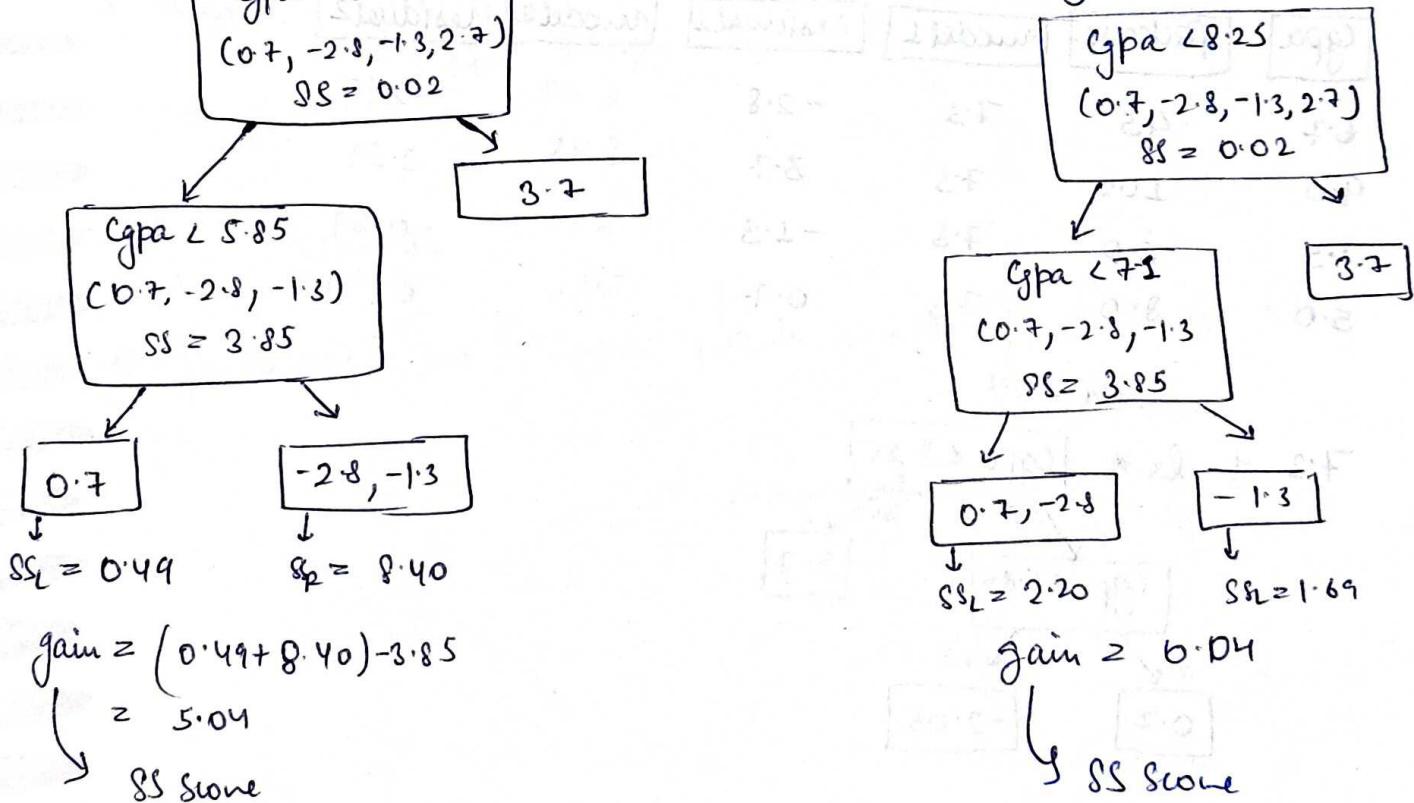
$$= 5.04$$

SS score

accuracy



② Splitting Criteria  $\rightarrow 7.1$

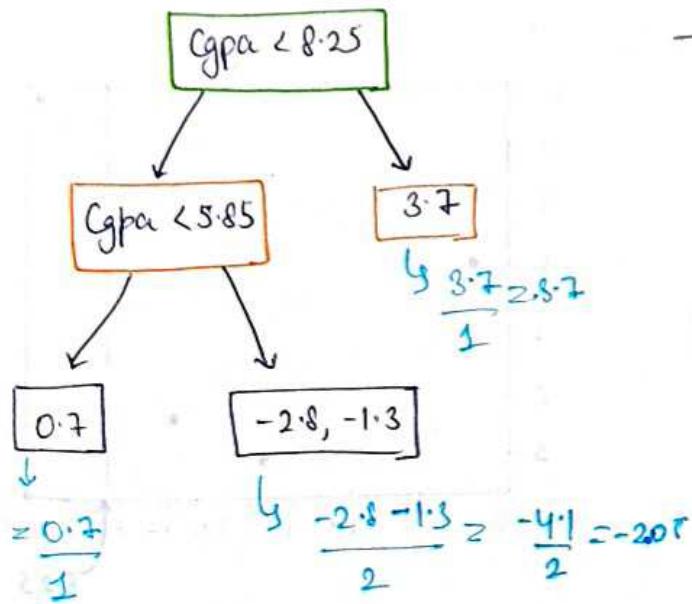


$$\text{gain} = 6.04$$

SS score

accuracy

# The final Decision Tree

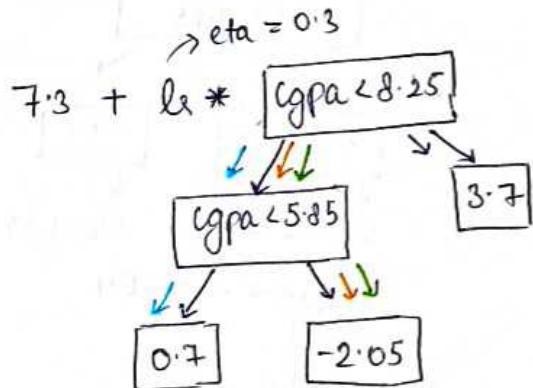


→ Now, we have to find out the output of all leaf nodes.

$$\text{Output} = \frac{\text{sum of residual}}{\text{no. of residual} + \lambda}$$

let  $\lambda = 0$

cgpa	package	model1	residual1	model2	residual2
6.7	4.5	7.3	-2.8	6.69	-2.19
9.0	11.0	7.3	3.7	8.41	2.59
7.5	6.0	7.3	-1.3	6.69	-0.69
5.0	8.0	7.3	0.7	7.51	0.49



$$7.3 + 0.3 * (-2.05) = 6.69$$

$$7.3 + 0.3 * (8.41) = 8.41$$

$$7.3 + 0.3 * (-2.05) = 6.69$$

$$7.3 + 0.3 * (0.7) = 7.51$$

### Stage ③

cgpa | res<sup>2</sup>  
 ↳ decision tree

→ if residual is going to near 0, then our dt is good

$$\text{mean} + \text{var} * \text{dt}_1 + \text{var} * \text{dt}_2$$

no estimator → decide no. of tree

↳ Algo 1: Exact Greedy Algorithm for split finding  
 (Covered this method) for small data set

↳ Algo 2: Approximate Algorithm for split finding  
 (for large dataset)

### Research

cgpa | package → cgpa | f<sub>1</sub> | f<sub>2</sub> | 12<sup>th</sup> mark | package  
 or  
 categorical

# XGBoost for Classification

cgpa	placed
------	--------

5.70 0

6.25 1

7.10 0

8.15 1

9.60 1



$$\text{Model} + \underbrace{\mu_1}_{\text{(Mean)}} + \underbrace{\log \frac{n_1}{n_2}}_{\text{log odds}} + \dots + \underbrace{\log \frac{n_2}{n_3}}_{\text{log odds}}$$

Xgboost  $\rightarrow$  different decision tree { similarity of score }

## Step by Step Solution

cgpa	placed	fined 1	Stage 1:
------	--------	---------	----------

5.70 0 0.405 base estimator  $\rightarrow$  mean

6.25 1 0.405 log(odd)  $\rightarrow$   $\log\left(\frac{P}{1-P}\right)$  P  $\rightarrow$  prob +ve class

7.10 0 0.405

8.15 1 0.405  $= \log\left(\frac{3/5}{2/5}\right) = \log\left(\frac{3}{2}\right) = 0.405$

9.60 1 0.405  $= 0.405$

$\downarrow$  log(odd)

we have to change into probability

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

$$= \frac{e^{0.405}}{1 + e^{0.405}} = 0.60$$

gpa	placed	Pred 1	Pred 1 (prob)	res
-----	--------	--------	---------------	-----

5.70	0	0.405	0.6	-0.6
6.25	1	0.405	0.6	0.4
7.10	0	0.405	0.6	-0.8
8.15	1	0.405	0.6	0.4
9.60	1	0.405	0.6	0.4

Stage 1 → errors

Pseudo residual

+   $\xrightarrow{\text{decision}}$  Tree

Xgboost

gpa	res	Pred 1(prob)
-----	-----	--------------

5.70	-0.6	0.6
6.25	0.4	0.6
7.10	-0.6	0.6
8.15	0.4	0.6
9.60	0.4	0.6

leaf node

-0.6, 0.4, 0.6, 0.4, 0.4

$\Leftrightarrow$  similarity score

$$SS = \frac{(\sum \text{residual}_i)^2}{\sum \text{prev-prob}_i(1-\text{prev-prob}_i)}$$

$$= \frac{(-0.6+0.4+(-0.6)+0.4+0.4)^2}{5 \times [0.6 \cdot (1-0.6)]} = 0$$

$\Rightarrow$  similarity score

gpa	5.70	5.97
	6.25	6.67
	7.10	7.62
	8.10	8.87
	9.60	

Splitting criteria = 5.97

$\boxed{\text{gpa} < 5.97}$

$\boxed{-0.6}$

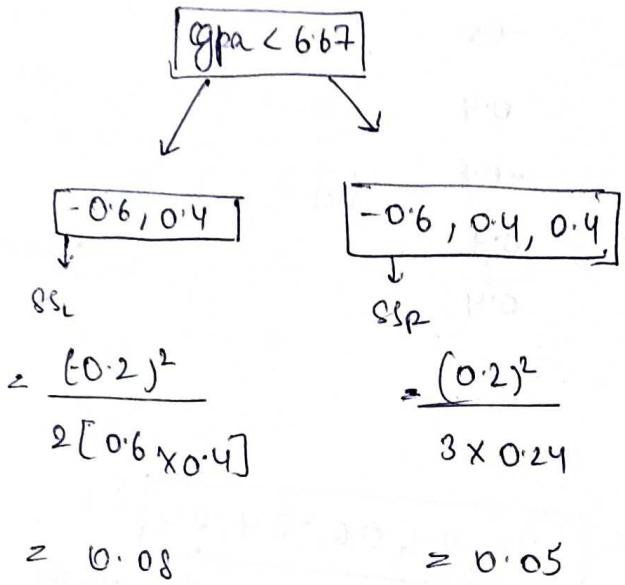
$$SS_L = \frac{(-0.6)^2}{0.6(1-0.6)} = 1.5$$

$\boxed{0.4, -0.6, 0.4, 0.4}$

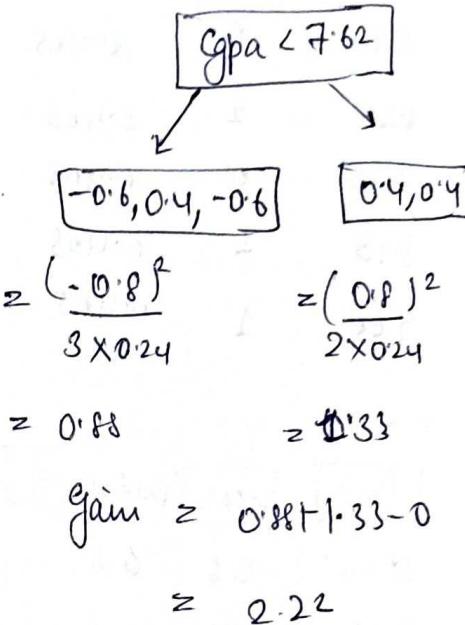
$$SS_R = \frac{0.4 - 0.6 + 0.4 + 0.4}{4 \times 0.6 \times 0.4} = 0.37$$

$$\text{gain} = (SS_L + SS_R) - SS_{\text{root}} \\ = 1.87$$

Splitting Criteria = 6.67

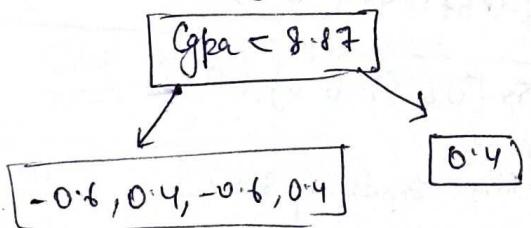


Splitting Criteria = 7.62



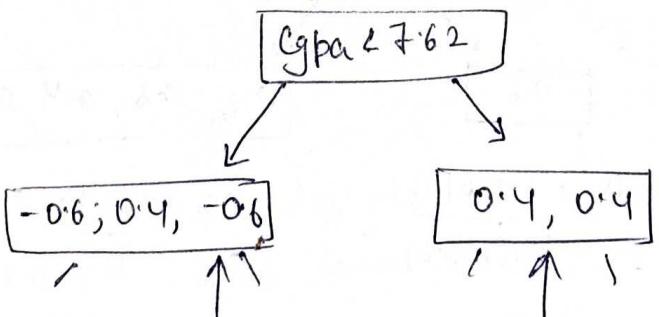
Criteria  
Splitting = 8.87

→ Same steps for this decision tree.



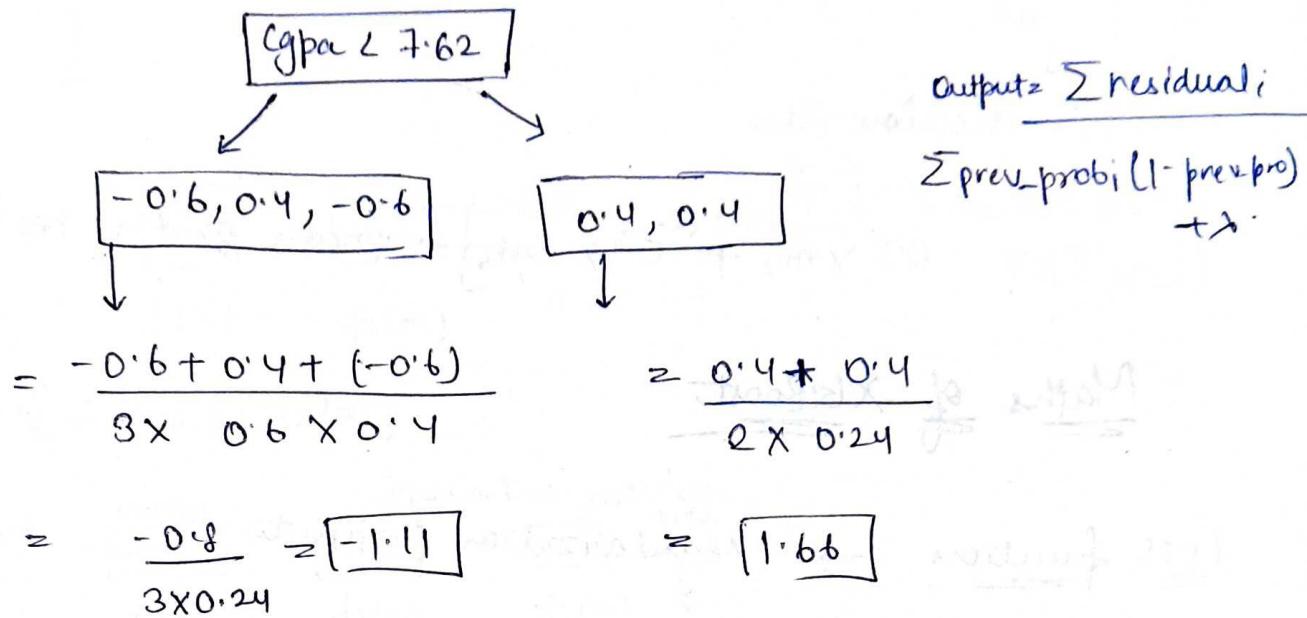
→ Splitting Criteria = 7.62 in a largest similar scene

Stage 2 model



\* We can also derive or find next node of these leaf node

## Stage 2 model



<u>cgpa</u>	<u>placed</u>	<u>Pred1 (log odds)</u>	<u>Pred1 (prob)</u>	<u>res1</u>	<u>Pred2 (log)</u>	<u>Pred2 (prob)</u>	<u>res2</u>
5.70	0	0.405	0.6	-0.6	0.072	0.518	-0.51
6.25	1	0.405	0.6	0.4	0.072	0.518	0.48
7.10	0	0.405	0.6	-0.4	0.072	0.518	-0.51
8.15	1	0.405	0.6	0.4	0.0923	0.712	0.28
9.60	1	0.405	0.6	0.4	0.903	0.712	0.48

$$\text{log (odds)} + 0.3 \times m_2 \rightarrow \text{learning rate}$$

$$0.405 + 0.3 \times \boxed{\text{gpa} < 7.62}$$

$$= 0.405 + 0.3 * (-1.11) \\ = 0.072$$

$$= 0.405 \cancel{+} 0.3 (-1.11) \\ = 0.032$$

$$= 0.405 + 0.3 \times 1.66 \\ = 0.903$$

$$\text{Prob} = \frac{e^{\text{log odds}}}{1 + e^{\text{log odds}}}$$

### Stage 3

Cgpa | res2

↳ decision Tree

$$[m_1 + 0.3 \times m_2 + 0.3 \times m_3] \rightarrow \text{pred} \rightarrow \text{pred} \rightarrow \text{res3}$$

### Maths of XGBoost

Loss function → regularization parameter

XGBoost loss function different from Gradient Boost

minimize

$m_1 + m_2$   
↓      ↓  
base      decision  
tree

$$L = \sum_{i=1}^n L(Y_i, \hat{Y}_i) + \boxed{\gamma L(\hat{t}_k(x_i))}$$

$L(A) = Y T + \frac{1}{2} \lambda \|w\|^2$   
 regularization parameter  
 no. of leaf node  
 leaf weight or output  
 regularization parameter

### The Derivation

$$L = \boxed{\sum_{i=1}^n L(Y_i, \hat{Y}_i)} + \boxed{\gamma L(\hat{t}_k)}$$

↓  
 loss function      ↓  
 objective function  
 regularization

mean  
 $m_1$   
 $f_1(x)$

$$\xrightarrow{\text{loss function}} \sum_{i=1}^n L(y_i, \hat{y}_i) \quad [\because \text{we have only mean not any decision tree}]$$

Case 1  $\rightarrow$

$$\underline{\text{Case 2}} \rightarrow \begin{array}{c} \text{mean} \\ m_1 \\ f_1(x_i) \end{array} \quad \begin{array}{c} \text{mean} \\ m_2 \\ \xrightarrow{\text{decision tree}} f_2(x_i) \end{array} \rightarrow L = \sum_{i=1}^n L(y_i, f_1(x_i) + f_2(x_i)) + R(f_2(x_i))$$

$$\hat{y}_i = f_1(x_i) + f_2(x_i)$$

$$\underline{\text{Case 3}} \rightarrow \begin{array}{c} \text{mean} \\ m_1 \\ f_1(x_i) \end{array} \quad \begin{array}{c} \xrightarrow{\text{decision tree}} m_2 \\ f_2(x_i) \end{array} \quad \begin{array}{c} \xrightarrow{\text{decision tree}} m_3 \\ f_3(x_i) \end{array} \rightarrow L = \sum_{i=1}^n (y_i, f_1(x_i) + f_2(x_i) + f_3(x_i)) + R(f_3(x_i))$$

Case 4:-

$$\begin{array}{c} \text{mean} \\ m_1 \\ f_1(x_i) \end{array} \quad \begin{array}{c} \xrightarrow{\text{decision tree}} m_2 \\ f_2(x_i) \end{array} \quad \cdots \quad \begin{array}{c} \xrightarrow{\text{decision tree}} m_t \\ f_t(x_i) \end{array} \rightarrow L = \sum_{i=1}^n y_i, f_1(x_i) + f_2(x_i) \dots + f_t(x_i) + R(f_t(x_i))$$

$\downarrow \text{simplify}$

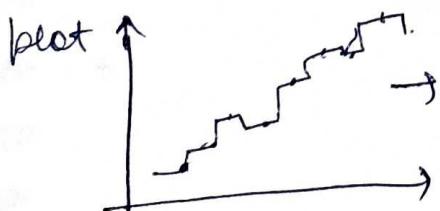
$$L^{(t)} = \sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)} + f_t(x_i) + R(f_t))$$

At any stage

$$L^{(t)} = \sum_{i=1}^n L(y_i, (\hat{y}_i^{(t-1)} + f_t(x_i))) + R(f_t(x_i))$$

$\rightarrow$  we have to the  $\hat{y}_i^{(t)}$  of  $f_t(x_i)$  which reduce value of  $L^{(t)}$ .

but it is difficult to reduce the value because in XGBoost we use decision tree and decision tree



non-differentiable

\* We have to smooth the line

Solution of the problem → we use Taylor Series

complex function → approx → polynomial

Applying Taylor Series

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

$$L \stackrel{\leftrightarrow}{=} \left[ \sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) \right] + \lambda (f_t(x_i))$$

$$x = \hat{y}_i^{(t-1)} + f_t(x_i) \quad a \rightarrow \hat{y}_i^{(t-1)}$$

$$f(a) = \sum_{i=1}^n L(y_i, \underbrace{\hat{y}_i^{(t-1)} + f_t(x_i)}_n)$$

↳ replace with a

$$f(a) = \sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)})$$

Gradient ( $\mathbf{g}_i$ )

$$f'(a)(x-a) = \sum_{i=1}^n \left[ \frac{\partial L(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}} \right] f_t(x_i)$$

$$\frac{f''(a)}{2} (x-a)^2 = \sum_{i=1}^n \left[ \frac{\partial^2 L(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)} \partial x_i} \right] f_t^2(x_i)$$

↳ hessian

$$L^{(t)} = \sum_{i=1}^m \left[ L(y_i, \hat{y}_i^{(t-1)}) + g_{ft}(x_i) + \frac{1}{2} h_{ft}^2(x_i) \right] + \Delta(f_t(x_i))$$

$$L^{(t)} = \sum_{i=1}^m \left[ g_{ft}(x_i) + \frac{1}{2} h_{ft}^2(x_i) \right] + \underbrace{\Delta(f_t(x_i))}_{\text{is expand by leaf node}}$$

$$L^{(t)} = \sum_{i=1}^n \left[ g_{ft}(x_i) + \frac{1}{2} h_{ft}^2(x_i) \right] + \gamma T + \frac{1}{2} \lambda [w_1^2 + w_2^2 + w_3^2 + \dots + w_T^2]$$

$$L^{(t)} = \sum_{i=1}^n \left[ \underbrace{g_{ft}(x_i) + \frac{1}{2} h_{ft}^2(x_i)}_{i=i+n} \right] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

So, there are 2 submission. We have make a common submission.

$$L^{(t)} = \sum_{j=1}^T \left[ \sum_{i \in I_j} g_i w_j + \frac{1}{2} \sum_{i \in I_j} h_i w_j^2 \right] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

\* Explaining How we make a common submission

$$L^{(t)} = \sum_{i=1}^n \left[ g_{ft}(x_i) + \frac{1}{2} h_{ft}^2(x_i) \right] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

$$w_j = f_t(x_i)$$

$$\sum_{i=1}^n \rightarrow \sum_{j=1}^T \sum_{i \in I_j}$$

$$L^{(t)} = \sum_{j=1}^T \left[ \sum_{i \in I_j} g_i w_j + \frac{1}{2} \sum_{i \in I_j} h_i w_j^2 \right] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

expand all now one by one

$$\rightarrow (g_1 f_t(x_1) + \frac{1}{2} h_1 f_t^2(x_1)) + (g_2 f_t(x_2) + \frac{1}{2} h_2 f_t^2(x_2)) + \dots + (g_m f_t(x_m) + \frac{1}{2} h_m f_t^2(x_m))$$

expand

row 1

row 2



Low  $\rightarrow$  acc. to GPA

$$w_1 = 5$$

$$w_2 = 10$$

① ② ③

$j=1, j=2$

$$\frac{g_{\text{soft}}(x_3) + \frac{1}{2} h_3 f^t(x_3)}{\downarrow \text{row 3}} + \frac{g_{\text{soft}}(x_1) + \frac{1}{2} h_1 f^t(x_1)}{\downarrow \text{row 1}} + \dots - \frac{g_{\text{soft}}(x_4) + \frac{1}{2} h_4 f^t(x_4)}{\downarrow \text{row 4}}$$

In this method, we can also cover all.

$f^t(x_i) \rightarrow$  Output of decision Tree ( $w_i$ )

So, we can write

$$\frac{g_{\text{soft}}(x_3) + \frac{1}{2} h_3 w_1^2}{\downarrow \text{row 3}} + \frac{g_{\text{soft}}(x_1) + \frac{1}{2} h_1 \cancel{w_1^2} (w_2)^2}{\downarrow \text{row 1}} + \dots - \frac{}{\downarrow \text{row 4}}$$

$$L^{(t)} = \sum_{j=1}^T \left[ \sum_{i \in T_j} g_i w_i + \frac{1}{2} \sum_{i \in S_j} h_i w_i^2 \right] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

We have to simplify this  $\uparrow$

$$L^{(t)} = \sum_{j=1}^T \left[ \sum_{i \in T_j} g_i w_i^2 + \frac{1}{2} \sum_{i \in S_j} h_i w_i^2 + \frac{1}{2} \lambda w_j^2 \right] + \gamma T$$

$$L^{(t)} = \sum_{j=1}^T \left[ \sum_{i \in T_j} g_i w_i + \frac{1}{2} \left( \sum_{i \in T_j} h_i + \lambda \right) w_j^2 \right] + \gamma T$$

$$\frac{\partial L^{(t)}}{\partial w_j} = \sum_{i \in S_j} g_i + \left( \sum_{i \in T_j} h_i + \lambda \right) w_j = 0$$

$$w_j = \boxed{\frac{- \sum_{i \in S_j} g_i}{\sum_{i \in T_j} h_i + \lambda}}$$

## Output value for Regression

$$w_j = \frac{-\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}$$

↳ residual

$$g_i = \frac{\partial L(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}}$$

$$\lambda = \frac{\partial^2 L(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)} \partial \hat{y}_j^{(t-1)}}$$

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

↳  
Reg. Loss  
Loss  
Mse

$$g_i = \frac{\partial L}{\partial \hat{y}_i} = -\sum_{i \neq j} (y_i - \hat{y}_i) = \sum_{i \neq j} (\hat{y}_i - y_i)$$

$$g_i = \sum_{i \in I_j} (y_i - \hat{y}_i^{(t-1)}) \rightarrow \text{sum of residual}$$

↳ residual  
for every point

$$h_{ij} = \frac{\partial L}{\partial \hat{y}_j} \Rightarrow \frac{\partial L}{\partial \hat{y}_j} = \sum_{i \neq j} (\hat{y}_i - y_i)$$

$$w_j = \frac{\sum_{i \in I_j} (y_i - \hat{y}_i^{(t-1)})}{\sum_{i \in I_j} 1 + \lambda} \Rightarrow \frac{\text{sum of residual}}{\text{no. of residual} + \lambda}$$

if no. of residual  
2 then 2

$w_j = \frac{\text{sum of residual}}{\text{no. of residual} + \lambda}$

for  
Regression

## Output value for Classification

$$w_j = - \frac{\sum_{i \in I_j} g_i}{\sum w_i + \lambda}$$

$$g = \frac{\partial L(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}}$$

$$h = \frac{\partial^2 L(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)} \partial \hat{y}_i^{(t-1)}}$$

loss - log loss  $\rightarrow$  gradient boosting

$$w_j = \frac{\text{Sum of residual}}{\sum_{i \in I_j} p_i(1-p_i) + \lambda}$$

$p_i \rightarrow$  Pred Prob of previous timestep

## Derivation of Similarity Score

$$L^{(t)} = \sum_{j=1}^T \left[ \left( \sum_{i \in I_j} g_i \right) w_j + \frac{1}{2} \left( \sum_{i \in I_j} w_i + \lambda \right) w_j^2 \right] + \chi T$$

$$w_j = \frac{-\sum_{i \in I_j} g_i}{\sum w_i + \lambda}$$

$$L^{(t)}(q_j) = \sum_{j=1}^T \left[ \left( \sum_{i \in I_j} g_i \right) \left( -\frac{\sum_{i \in I_j} g_i}{\sum w_i + \lambda} \right) + \frac{1}{2} \left( \sum_{i \in I_j} w_i + \lambda \right) \left( \frac{\sum_{i \in I_j} g_i}{\sum w_i + \lambda} \right)^2 \right]$$

$$\frac{1}{2} \left( \sum_{i \in I_j} w_i + \lambda \right) \left( \frac{\sum_{i \in I_j} g_i}{\sum w_i + \lambda} \right)^2 + \chi T$$

$$\chi^{(t+1)}(q_j) = \sum_{j=1}^T \left[ -\frac{\left( \sum_{i \in I_j} g_i \right)^2}{\sum w_i + \lambda} + \frac{1}{2} \frac{\left( \sum_{i \in I_j} g_i \right)^2}{\sum w_i + \lambda} \right] + \chi T$$

$$= \sum_{j=1}^T \left[ -\frac{1}{2} \frac{\left( \sum_{i \in I_j} g_i \right)^2}{\sum_{i \in I_j} h_i + \lambda} \right] + \gamma T$$

No. of  
nodes.  
parent / child

$$\mathcal{L}^{(t)}(a) = -\frac{1}{2} \sum_{j=1}^T \left[ \frac{\left( \sum_{i \in I_j} g_i \right)^2}{\sum_{i \in I_j} h_i + \lambda} \right] + \gamma T$$

we can replace g with entropy

neg

$g \rightarrow$  sum of residual

$n \rightarrow \sum n \rightarrow$  No. of residual

classi

$\downarrow$   
 $g \rightarrow$  residual

$$h = p_i (1-p_i)$$