

Lecture - 05 :

→ Ensemble Techniques [Bagging, Boosting]

→ Random forest.

→ AdaBoost

→ Xg Boost.

Ensemble Techniques:

→ Classification and Regression Problem.

We are Solved One algorithm for any Classification and Regression Problem.

Can we Solve multiple algorithm to solve a problem?

Yes, we can use Multiple algorithm to solve a problem. Something called as Ensemble methods.

Ensemble methods



The Two techniques used to solve problems with multiple algorithms.

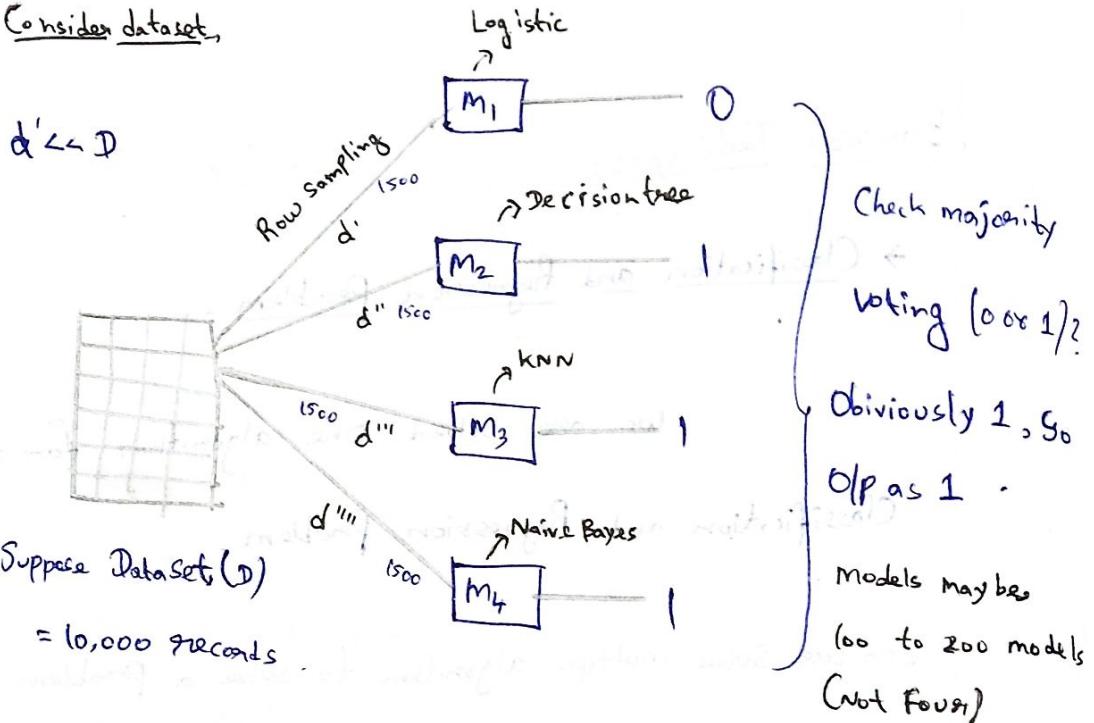
BAGGING:

(Faster & better) equivalent statement
Bagging is a technique to combine multiple algorithm

to solve a problem. [parallelly]

Classification
Problem

Consider dataset,

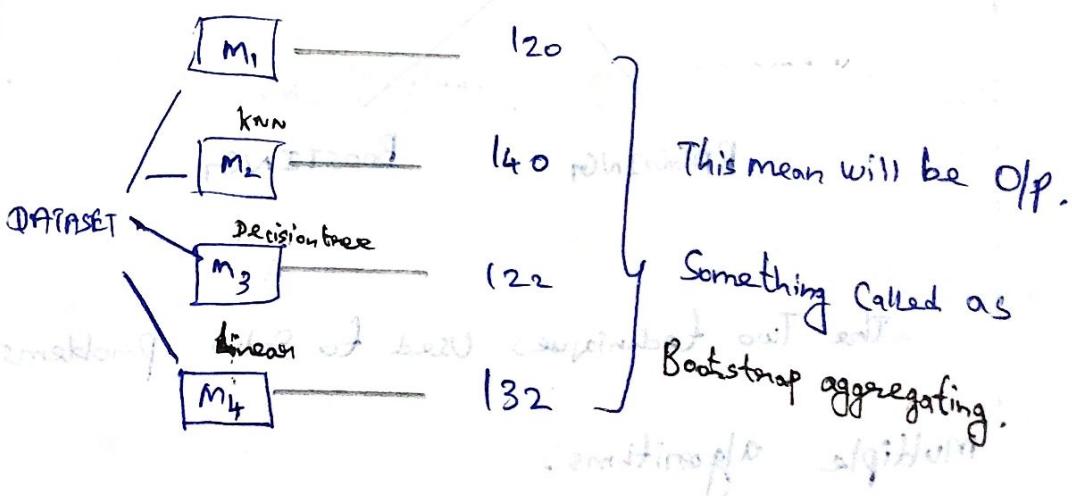


The records in d' or d'' or d''' may be in d'' or may be in d''' !

Just row sampling.

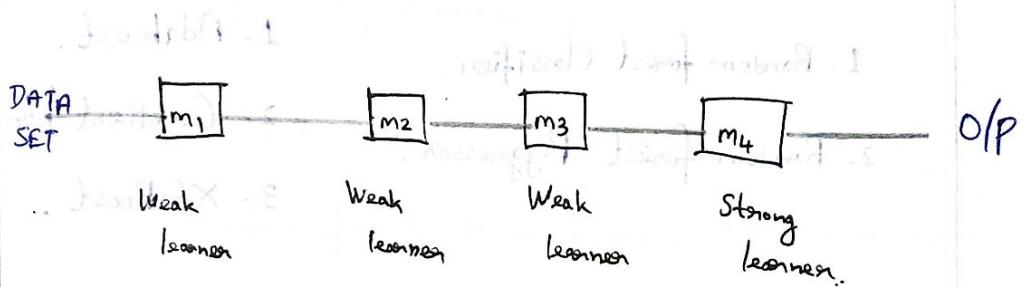
The above thing is for the Classification problem.

Linear



BOOSTING

Bagging is Combined model as Parallelly, Then the Boosting Combined Model Serially and also improve the efficiency of output.



To Combine these weak learners to get the Strong learner and get efficient Output.

Example

To create a data Science Course, we need



Another Example, One Scientist Solve a particular problem which

involves Various geography, maths, Physics and Chemistry.



Many Experts are help to find the answer

And give Correct answer.

Which algorithm is used for working with BAGGING

and BOOSTING?

What are the main differences between them?

BAGGING



1. Random forest classifier.
2. Random forest Regressor.

BOOSTING



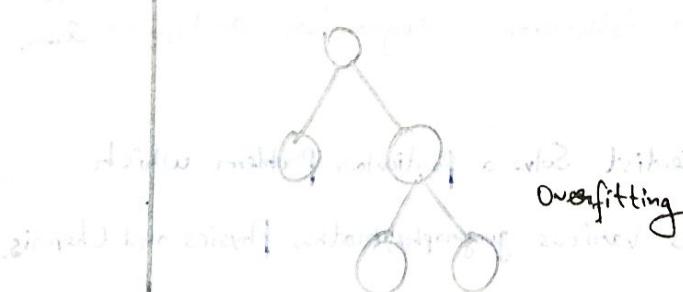
1. Adaboost.
2. Gradient boost.
3. XGBoost.

BAGGING,

Random forest Classifier and Regressor:

Before moving, what is the main problem in the

Decision Tree?



So, Generally, the overfitting
is occurred.

In order to avoid overfitting,

We use Pruning techniques.

The Post Pruning is used,

Sometime the large dataset.

We Can't Post Pruning.

On Otherhand, Pre Pruning is
Set by hyperparameter.

Sometime the model may
be not accurate based
on this data.

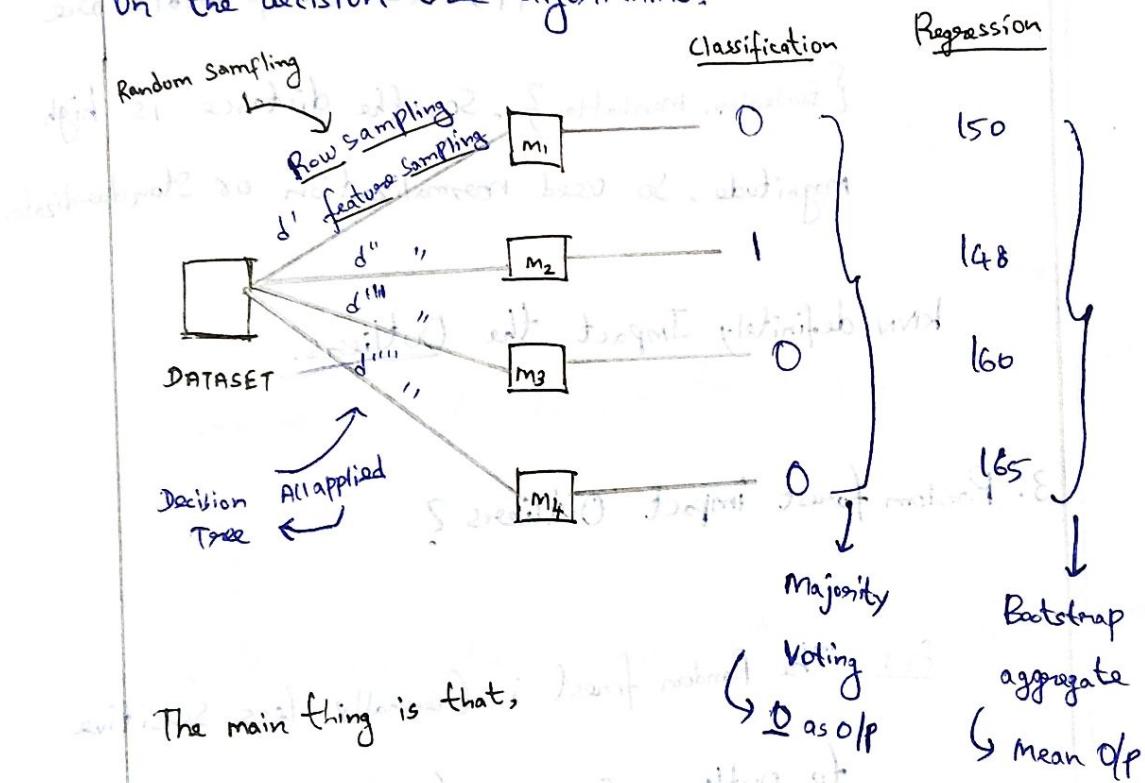
In this case, the random forest helps to produce the High Variance to Low Variance.

Then we get the Generalised Model.



How works? Work?

The random forest is a bagging technique is applied on the decision tree algorithms.



That's why we called it Random forest.

→ This helps to avoid Overfitting in

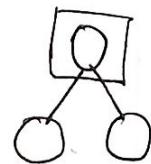
the decision tree and get Generalised model.

Some Interview Questions:

1. Is Normalization required for Decision Tree?

Ans: No, because the decision tree split the things

Only. [Not apply any Quantity].



2. Is kNN is required Standardization or Normalization?

Ans: Yes, Because we apply the two formulae are

{ Euclidean, Manhattan }. So, if the distance is high

Magnitude. So used normalization or Standardization.

kNN definitely Impact the Outliers.

3. Random forest impact Outliers?

Ans: The Random forest is Generally less sensitive

to outliers compared to individual decision

tree.

But Decision Tree is impacted by outliers due to

the nature of making Splits.

In Bagging, Custom Bagging is also available.

BOOSTING:

with weak learner and weak

ADABOOST:

Adaboost also working with Decision Tree.

Let's consider,

| f_1 | f_2 | f_3 | f_4 | O/P | Weight |
|-------|-------|-------|-------|-----|---------------|
| - | - | - | - | Yes | $\frac{1}{7}$ |
| - | - | - | - | No | $\frac{1}{7}$ |
| - | - | - | - | - | $\frac{1}{7}$ |
| - | - | - | - | - | $\frac{1}{7}$ |
| - | - | - | - | - | $\frac{1}{7}$ |
| - | - | - | - | - | $\frac{1}{7}$ |
| - | - | - | - | - | $\frac{1}{7}$ |

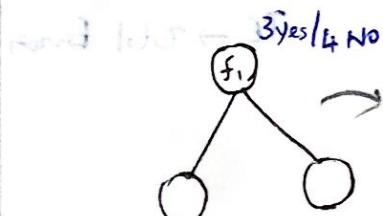
Total records = 7 [Based on this weights are assigned]

$$\text{Overall weight} = 1 \quad \left(\frac{1}{7} + \frac{1}{7} + \frac{1}{7} + \frac{1}{7} + \frac{1}{7} + \frac{1}{7} + \frac{1}{7} = 1 \right)$$

After that taking the feature with the help of information

Gain.

In this Scenario, the decision tree take,



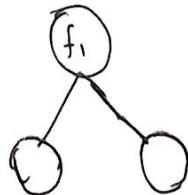
In Adaboost, The decision tree

Only take one depth of the tree. This is called as

STUMPS.

Stump - One level decision Tree.

Why we Stump? Because it is weak learner.



| f_1 | f_2 | f_3 | f_4 | O/P | Weight | New weight |
|-------|-------|-------|-------|-------|---------------|------------|
| - | - | - | - | yes | $\frac{1}{7}$ | 0.05 |
| - | - | - | - | No | $\frac{1}{7}$ | 0.05 |
| - | - | - | - | - | $\frac{1}{7}$ | 0.05 |
| - | - | - | - | No | $\frac{1}{7}$ | 0.349 |
| - | - | - | - | wrong | $\frac{1}{7}$ | 0.05 |
| - | - | - | - | - | $\frac{1}{7}$ | 0.05 |
| - | - | - | - | - | $\frac{1}{7}$ | 0.05 |

After train this weak learner

Model, Pass all the records

Suppose one of the record

Predict wrong.

The 4th row O/P is No, Suppose predicting "yes" as a wrong.

If wrong, Calculating the Total Error, ①

1. Calculating Total Error, 1 = if find wrong

$$\text{Total Error} = \frac{1}{7} \cdot [weight] \quad \text{point that wrong}$$

2. Performance of Stump,

$$\text{error rate} = \frac{1}{2} + \log_e \left(\frac{1 - TE}{TE} \right) \quad \text{where } TE \rightarrow \text{Total Error}$$

$$= \frac{1}{2} + \log_e \left(\frac{1 - \frac{1}{7}}{\frac{1}{7}} \right)$$

$$= 0.895$$

3. Updating the Weights,

Why Updating? Because find out the correct and wrong records. If wrong record found, ~~skip~~ this pass to another weak learner. [Stumps]

$P_S \rightarrow$ Performance of STUMP.

Correct records:

$$\text{New sample weight} = \text{Weight} * e^{-P_S} = \frac{1}{7} * e^{-0.895} = 0.05$$

$$\frac{1}{7} = 0.1428 \downarrow \downarrow \text{reduced}$$

$$0.05$$

Incorrect records:

$$\text{New sample weight} = \text{Weight} * e^{P_S} = \frac{1}{7} * e^{0.895} = 0.349$$

New weight: See the left page.

New weight

Normalized weight

Buckets

0.05

$$0.05 / 0.649 = 0.077$$

[0 - 0.07]

0.05

$$0.05 / 0.649 = 0.077$$

[0.07 - 0.14]

0.349

$$0.349 / 0.649 = 0.537$$

[0.14 - 0.21] High size

0.05

$$0.077$$

[0.747 - 0.753]

0.05

$$0.077$$

[0.753 - 0.760]

0.05

$$0.077$$

[0.760 - 0.767]

0.649

$$\underline{\underline{\approx 1}}$$

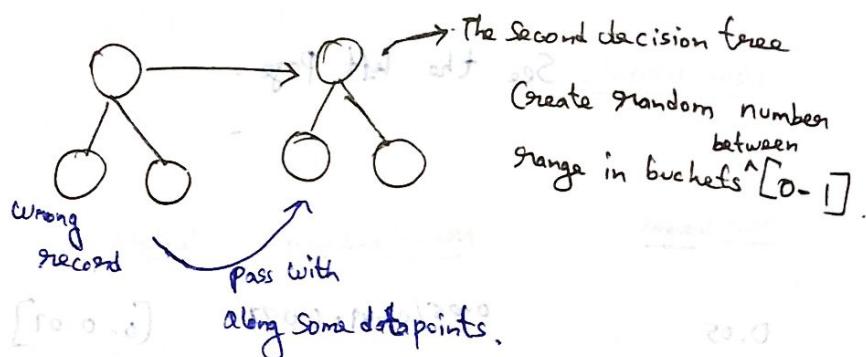
$\Sigma L = \text{error} \times \text{threshold}$

Why Normalize?

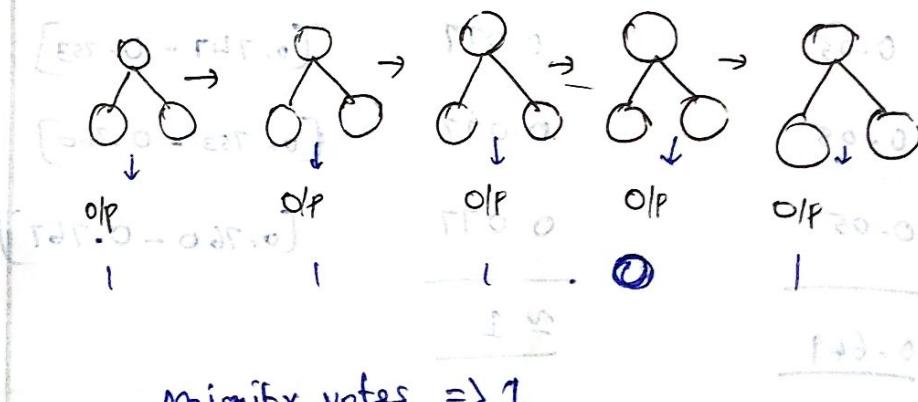
While update the weights, the magnitude ~~are~~ is ~~not~~ changed. So, the new weight does not come to 1. That's why normalize them.

Buckets are used to identify the wrong records, if conclude the bucket range is high.

4th row record size bucket is high. So, identified 4th record as wrong and error.



Then the second decision tree train with wrong records and also some datapoints. After some amount of decision tree, the output as,



So, 1 as the output for classification problem.

In regression,

O/P O/P O/P O/P O/P
66-3 69 68-3 67 59

Aggregate Mean



O/P for regression problem.

So, the number of decision tree, the time complexity is worst.

NOTE: BAGGING AND BOOSTING not only support random forest and decision tree and also support and algorithm for base model.

BLACK BOX MODEL VS WHITE BOX MODEL: [Interview Q/A]

BLACK BOX - We cannot see how algorithm work deeply.

WHITE BOX - We can understand how algorithm work deeply.

Linear Regression → white box [ε value change] ANN → Black Box
Random forest → Black box [No of decision Tree]
Decision Tree → White Box [we see decision tree]

[Unable to
see how many
neurons
present]

BOOSTING:

GRADIENT BOOSTING:

Consider the dataset,

| Exp | Degree | Salary | <u>base model prediction \hat{y}</u> | Residuals | | | |
|-----|---------|--------|---|----------------|----------------|----------------|----------------|
| | | | | R ₁ | R ₂ | R ₃ | R ₄ |
| 2 | BE | 50k | 75 | -25 | -23 | - | - |
| 3 | Masters | 70k | 75 | -5 | -3 | - | - |
| 4 | Masters | 80k | 75 | 3 | 3 | - | - |
| 5 | Phd | 100k | 75 | 25 | 20 | - | - |

Residual
decreasing

Step 1:

Base model O/P.

$$\text{Average mean} = \frac{50 + 70 + 80 + 100}{4} = 75 \Rightarrow \hat{y}$$

Step 2:

Compute residuals Error, Pseudo residuals.

- $\hat{y} \rightarrow \text{Base model Residual} \rightarrow R_1, R_2, R_3, R_4$
- $\hat{y} \rightarrow \text{Least squares error and know residual as self - real value}$

Step 3:

Construct decision tree $\{x_i, R_i\}$

$x_i \rightarrow \text{Independent feature}$

$R_i \rightarrow \text{Residuals of base model}$

Both are passed decision tree.

Then the decision tree work with inputs and gives the residuals as output (R_2).

suppose $f(x)$ have some value -

suppose take the first record,

$$\begin{array}{l} \text{residuals} = \hat{y} - R_2 \\ \text{base model} = f(x) = \sum_{i=1}^n w_i x_i + b \\ \text{so } \hat{y} = \sum_{i=1}^n w_i x_i + b + R_2 \end{array}$$

We pass the first record in base model,

$$\begin{aligned} \hat{y} &= \hat{y} - R_2 \\ &= 75 - 23 \Rightarrow 52 \quad [\text{It is overfitting because it is closely to } y] \end{aligned}$$

To prevent overfitting,

To add Learning rate (α)

$$\begin{aligned} &= 75 + \alpha R_2 \\ &= 75 + 0.1 (-23) \\ &= 75 - 2.3 = 73.7 \quad [\text{But this value is huge diff from } y] \end{aligned}$$

So, we add another decision tree,

Therefore, this decision tree computes next

residuals $\left[\begin{array}{c} (75) - 73.7 \\ (75) - 73.7 \end{array} \right]$ Then give the prediction.

$$f(x) = \hat{y} + R_2$$

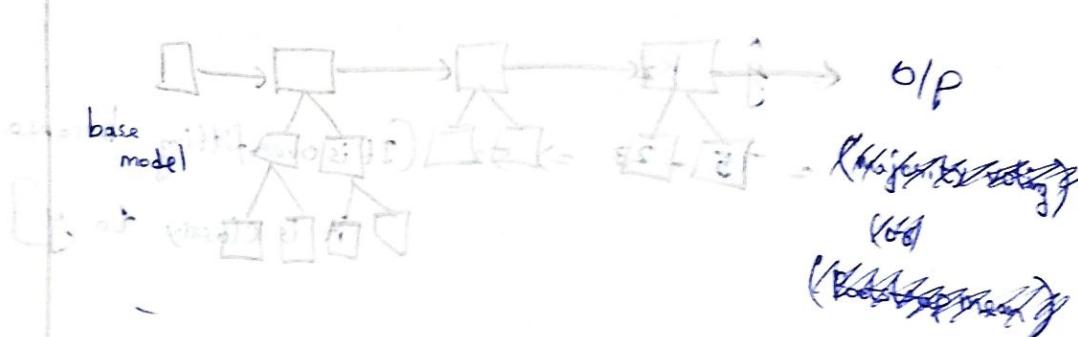
$$x = 3.5 \text{ not } 3.0$$

So, the formula would be,

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

↓ ↓ ↓ ↓ ↓
 base model model-1 model-2 ... Model-n
 {decision tree} {decision tree} {decision tree} {decision tree}

$$F(x) = \sum_{i=1}^n \alpha_i h_i(x)$$



Algorithm:

(*) start from initial value of f_0

1. Initialize model with a Constant value.

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

2. For $m=1$ to M :

① Compute so-called Pseudo-residuals :

first estimate and update it, repeat

$$\gamma_{im} = - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right] \text{ derivative}$$

$$f(x) = f_{m-1}(x)$$

for $i=1, 2, \dots, n$

2. $\{(x_i, y_{im})\}_{i=1}^n$ (Fit the base learner)

3. Compute multiplier γ_m

$$\gamma_m = \operatorname{argmin} \sum_{i=1}^n L(y_i, f_{m-1}(x_i) + \gamma h_m(x_i))$$

4. Update the model,

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

3. Output $f_m(x)$

Inputs:

1. $\{x_i, y_i\}$

$y_i \rightarrow$ dependent feature.

2. $L(y, f(x))$

↳ Loss function

Regression - Mse, mae etc

Classification - Log loss etc.

3. Number of trees in gradient boosting.

These three inputs are required for gradient boosting.

Consider,

↑ base model

| Exp | Degree | Salary | \hat{y} | How aligned? |
|-----|--------|--------|-----------|--------------|
| 2 | BE | 50k | 60 | Bad fit |
| 3 | Phd | 70k | 60 | Bad fit |
| 4 | Master | 60k | 60 | Exact fit |

(i) initial + (ii), and (iii)

$\gamma \rightarrow$ predicted value (\hat{y})

$$f_0(x) = \underset{\gamma}{\operatorname{argmin}} \left(\sum_{i=1}^n L(y_i, \gamma) \right)$$

↳ reduce the error

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

① Initialize model with Constant value Using above formula,

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

→ All records value.

$$= \frac{1}{2} (50 - \hat{y})^2 + \frac{1}{2} (70 - \hat{y})^2 + \frac{1}{2} (60 - \hat{y})^2$$

Which value is assigned for Constant? Ans: \rightarrow

Simple derivative, partial derivative of loss function w.r.t. Just derivative

$$\frac{\partial}{\partial x} (x^n) = nx^{n-1}$$

↓ first order

greatest

$$= \frac{2}{2} (50 - \hat{y}) (-1) + \frac{2}{2} (70 - \hat{y}) (-1) + \frac{2}{2} (60 - \hat{y}) (-1)$$

$$= -50 + \hat{y} - 70 + \hat{y} - 60 + \hat{y}$$

$$= -180 + 3\hat{y}$$

$$3\hat{y} = 180 \Rightarrow \hat{y} = 60$$

$$\boxed{\hat{y} = 60}$$

This value assigned for Constant base model. (\hat{y}) .

② Iterate the tree Count $M = 1$ to M .

↳ Compute Pseudo residuals ,

$$(x_i^*) = - \left[\frac{\partial L(y, f(x_i))}{\partial f(x_i)} \right] \text{ for } 1 \text{ to } n$$

It represents the

$$\text{Loss } \frac{1}{2} (y - \hat{y})^2$$

$$= \frac{2}{2} (y - \hat{y}) (-1)$$

Looks like

$$= -y + \hat{y}$$

$$\text{substituting value of } \frac{\partial L(y, \hat{y})}{\partial \hat{y}} = - (y - \hat{y})$$

$$\boxed{- \frac{\partial L(y, \hat{y})}{\partial \hat{y}}} = y - \hat{y}$$

(i) Basically, $\gamma_{im} = \arg \min_{\gamma} L(\hat{y}_i - (\gamma + f_m(x_i)))^2$

$$\begin{array}{c} \gamma_{11} = -10 \\ \downarrow \\ \text{first model } (x) \\ \text{second model } (\text{base model}) \end{array}$$

$$\gamma_{11} = -10$$

↳ 1st second.

$$\gamma_{21} = 10$$

$$\gamma_{21} = 10$$

↳ 2nd second

$$\gamma_{31} = 0$$

Once get the residual free, then fit base learner $f_m(x)$,

(ii) Learn and

↳ fit a Base learner $f_m(x)$

At this point not yet start

$$\text{if } \{(x_i, y_{im})\}$$

start from 0 degree

$$\textcircled{3} \quad \gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i - (\gamma + f_m(x_i)))^2$$

↳ Previous model

$$\text{Loss } \frac{1}{2} (y - \hat{y})^2$$

$$\sum_{i=1}^n \frac{1}{2} (y_i - (60 + \hat{y}))^2$$

↓
previous model output.

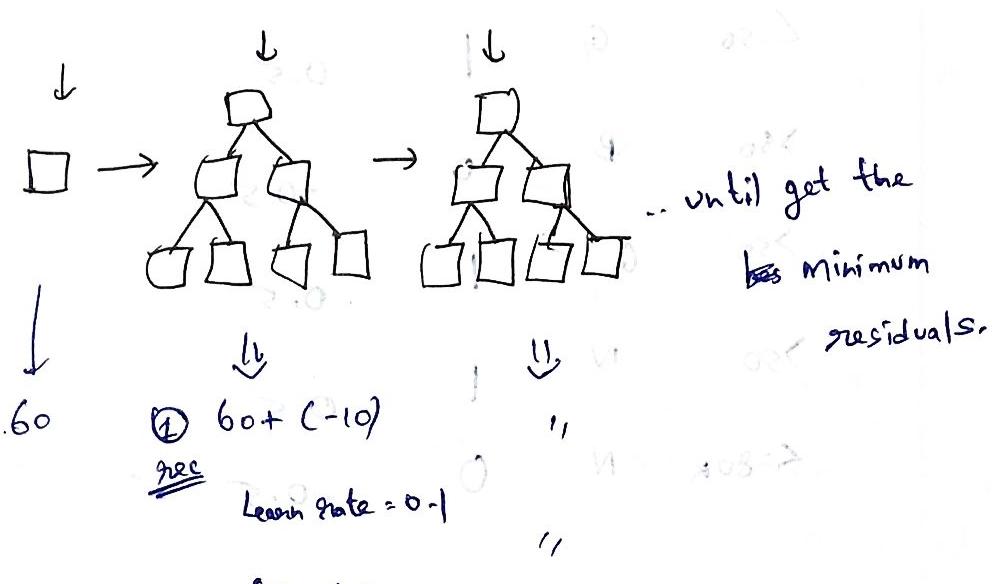
Then again calculate residues

$\arg \min \rightarrow$ It means minimum
 $B - B =$ residuals.

④ Update the model,

$$f_m(x) = f_{m-1}(x) + d(h(x))$$

Without $d\theta$, It may cause overfitting.



$$60 - 1.0$$

$$= 59.0$$

$$60 + (0.1)(10)$$

$$= 61$$

These all are about the Gradient Boosting.

BOOSTING,

XGBOOST:

xgboost Classifier:

Consider the dataset,

and which you fit

| <u>Salary</u> | <u>Credit</u> | <u>Approval</u> | <u>Residual</u> (R_1) | $\sqrt{(y - \hat{y})}$ | $P_S = 0.5$ | R_2 | R_3 | R_4 |
|---------------|---------------|-----------------|---------------------------|------------------------|-------------|-------|-------|-------|
| $<= 50k$ | B | 0 | 0.5 | 0.5 | | | | |

$L = 50k$ G ad node 0.5 (weak)

$L = 50$ G | 0.5

> 50 B 0 -0.5

> 50 G | 0.5

> 50 N B | 0.5

$<= 80k$ N 0 -0.5

Step 1: Create base model and given $O/P = 0.5$ (weak learner)

After creating base model, the following steps to be construct,

1. Create a Binary Decision Tree Using the feature.

2. Calculate Similarity weight,

$$SW = \sum (\text{Residual})^2$$

$$\sum (P_S(1-P_S) + \lambda)$$

3. Information Gain.

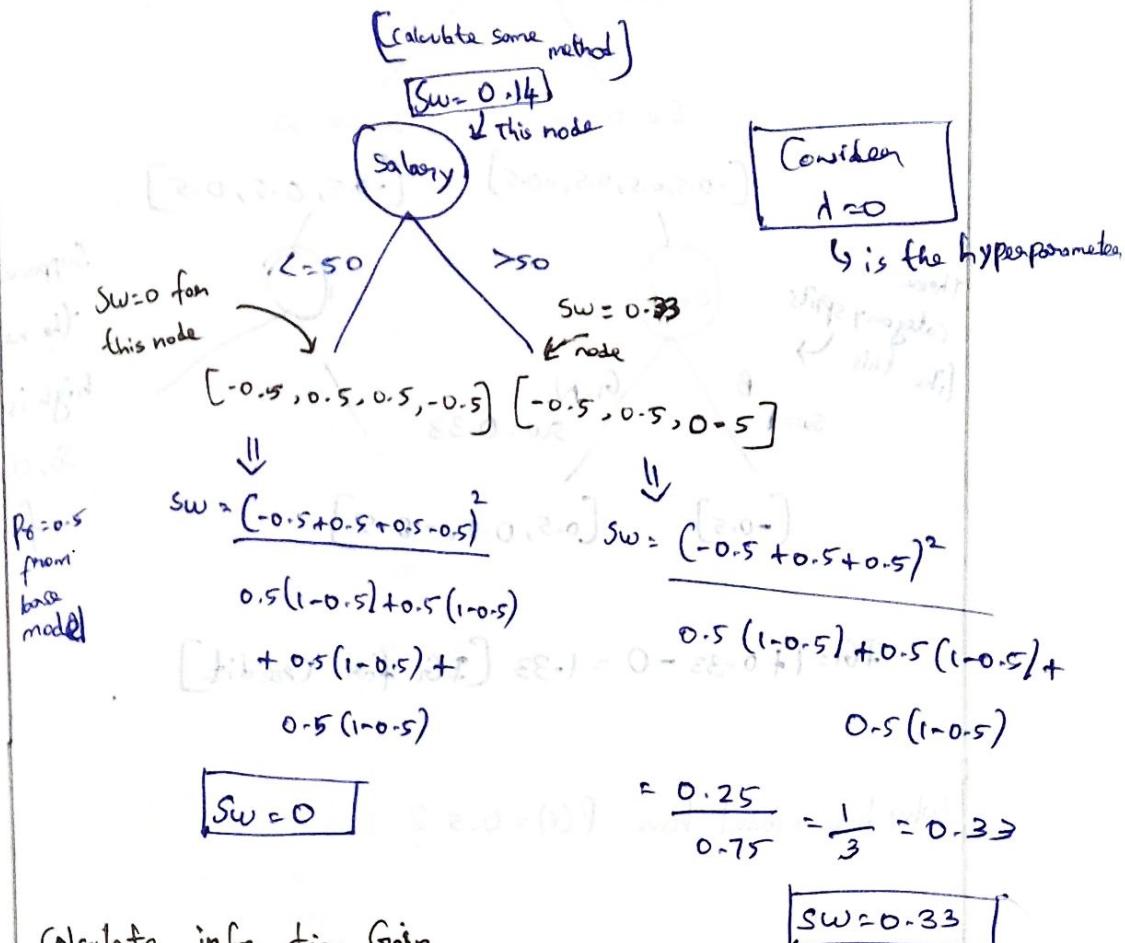
1. Base model - $O/P = 0.5$

$$P_S = 0.5$$

Ensure base model trained with the residuals before creating binary decision tree.

2. Create a Binary decision tree.

Consider Salary has a high information Gain.



Calculate information Gain,

This information finding is little different, we find information finding through Similarity weight,

$$IG_r = \sum_{\text{Childnode}} \frac{\text{Similarity weight of Childnode}}{SW}$$

$$IG_r = (0 + 0.33) - 0.14$$

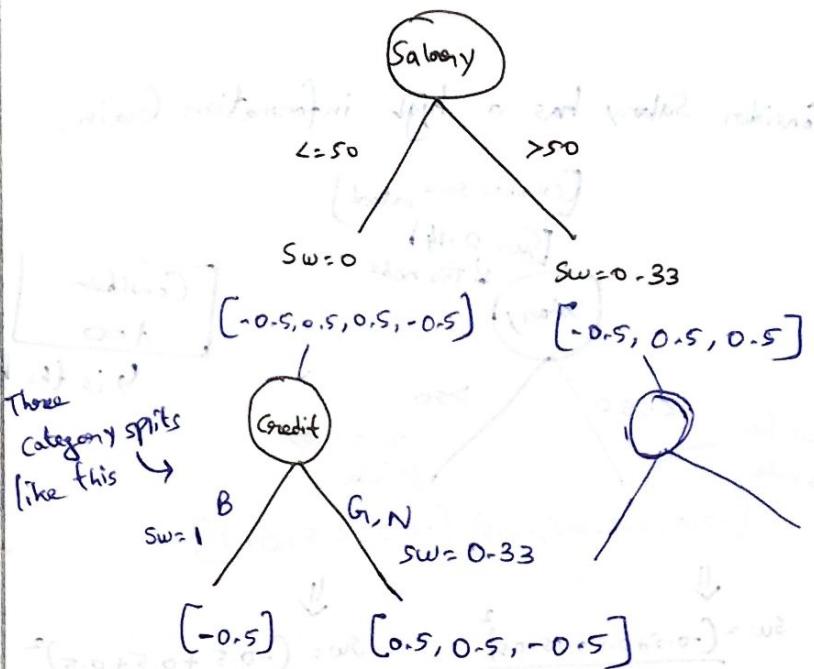
$$IG_r = 0.19 \rightarrow \text{It means this is the}$$

information Gain for root node.

[Salary]

Then again further split will be happened

and Continue.



Suppose Consider
the next IG_i is
high is Credit.
So, choose Credit
so for split.

$$IG_i = 1 + 0.33 - 0 = 1.33 \quad [\text{IG}_i \text{ for Credit}]$$

Why base model have $P(s) = 0.5$?

We Something apply, $\log\left(\frac{P}{1-P}\right) = \log\left(\frac{0.5}{0.5}\right) = 0$.

So, In first Case whenever the record goes to the base model, The model gives 0 with 0.5 probability.

So, The boosting look like,

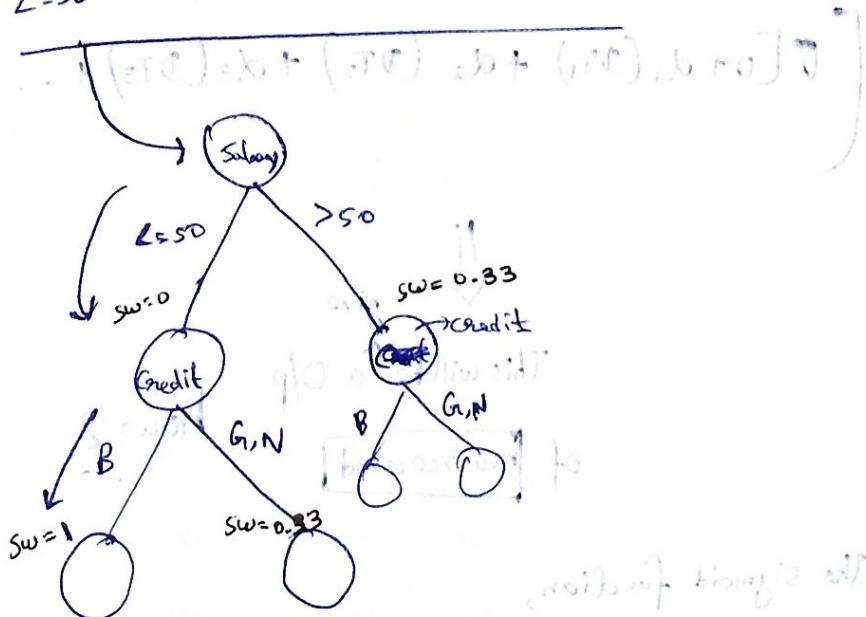
~~For each suppose the record~~

for understanding,

booster at this step next step will

Take the first record in the dataset.

| Salary | Credit | Approval | Residual |
|--------|--------|----------|----------|
| $L=50$ | B | 0 | -0.5 |



So, $[0 + \alpha(sw)]$

$[0 + \alpha(1)]$

base model first model
of decision tree

Then apply Sigmoid function, $\sigma[0 + \alpha(1)]$

The output comes between 0 to 1.

The decision trees are

more numbers. When stop?

The answer is Xgboost evaluate

the metrics and performance

based on that, The decision

tree making is stop at which node?

Then the model = decision tree as same,

$$\sigma[0 + \alpha(\text{DT}_1) + \alpha(\text{DT}_2)]$$

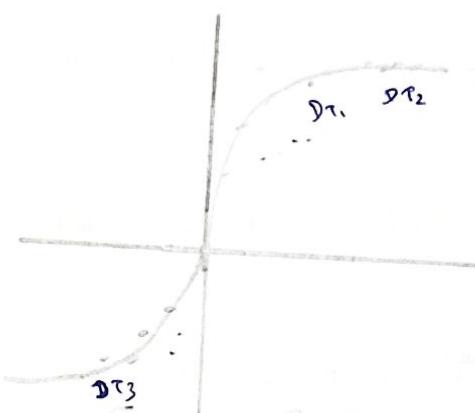
start

Likewise,

$$\left[\sigma[a + d_1(DT_1) + d_2(DT_2) + d_3(DT_3) + \dots + d_n(DT_n)] \right]$$

↓
give
This will be a O/p
of new record how?

The sigmoid function,



$$[(w_1)x_1 + w_2x_2 + b]$$

⇒ In this way, The O/p
will be classified.

$(w_1)x_1 + w_2x_2 + b = 0$ is not only linear if you not
This is how the Xgboost classifier work.

If it is the black box model, Because we didn't

See the working of decision tree

Pre Pruning → Applied Some hyperparameters tuning.

Left, (DT) feature has

Decision of left instance sit

stopping feature here.

noisy sit don't no best

Note:

The each decision tree is based on independent

feature not fake dependent One = /don't use

$$[(w_1)x_1 + w_2x_2 + b = 0]$$

XGBoost REGRESSOR: Base model o/p as average of all salary is 51.

Consider the dataset,

| Exp | Grp | Salary | Residual |
|-----|-----|--------|-------------------------|
| 2 | Yes | 40k | -11k (<u>40 - 51</u>) |
| 2.5 | Yes | 42k | -9k (<u>42 - 51</u>) |
| 3 | No | 52k | 1k (<u>52 - 51</u>) |
| 4 | No | 60k | 9k |
| 4.5 | Yes | 62k | 11k |

The similarity weight Change in regression, Consider,

$$\text{Similar weight } (sw) = \frac{\sum (\text{Residuals})^2}{\text{No of residuals} + \lambda}$$

Tree will be $= 25.011 + 25.881 = 50.89$

$$[-11, -9, 1, 9, 11]$$

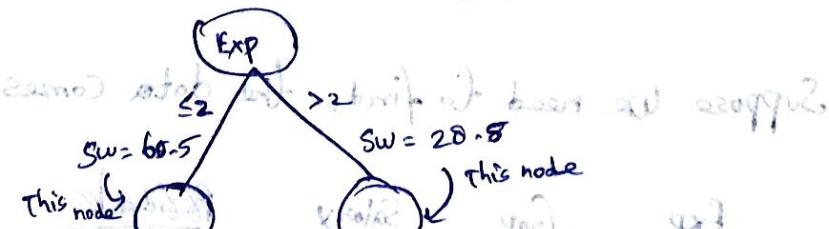
$$sw = 121$$

$$1+0$$

$$sw = 121$$

residuals \Rightarrow suppose $\lambda = 1$, $sw = 25.0142$, and no grid not

$$\text{Suppose take } \lambda = 1, sw = \frac{(-11 - 9 + 1 + 9 + 11)^2}{5+1} = \frac{16}{6} = 0.16$$



$$sw = \frac{121}{1+1}$$

$$sw = 60.5$$

$$sw = \frac{(-9 + 1 + 9 + 11)}{4+1} = \frac{144}{5} = 28.8$$

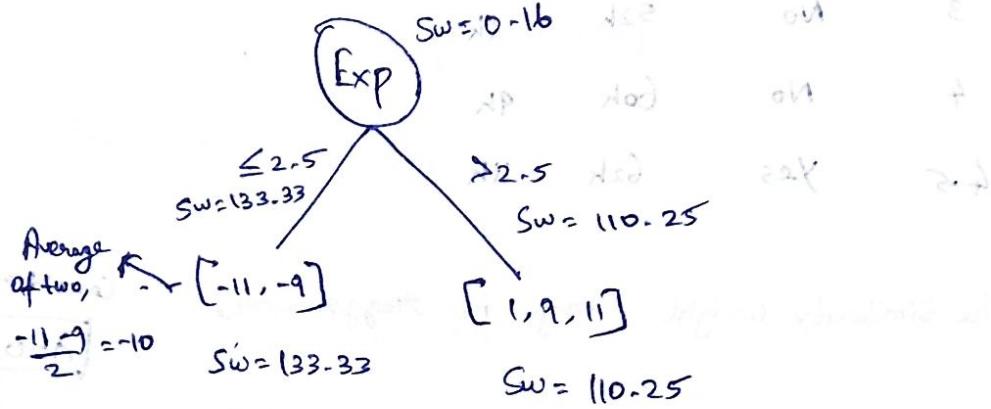
$$sw = 28.8$$

Then find the information gain,

$$= 65.5 + 28.8 - 0.16 = \underline{\underline{92.12}}$$

$$= \frac{89.13}{\underline{89.13}} \downarrow$$

Next split,



The information gain,

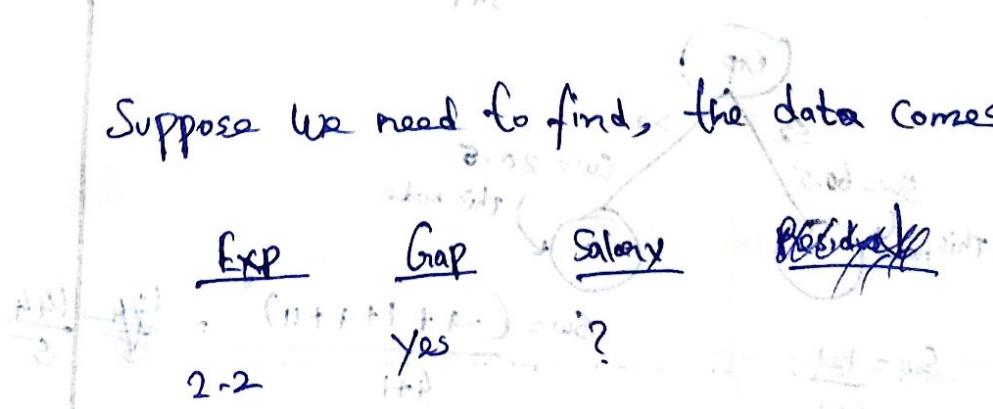
$$IG_1 = 133 - 33 + 110 - 25 = 0.16$$

IG = 243.42

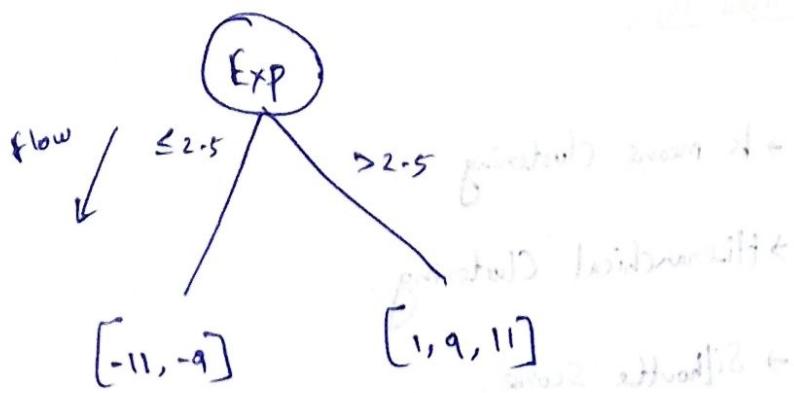
↳ This is greater IG compare to the previous IG.

for doing all the splits and possibilities whichever have a high information gain to take the split.

Suppose we need to find, the data comes



Suppose the above data comes the category,



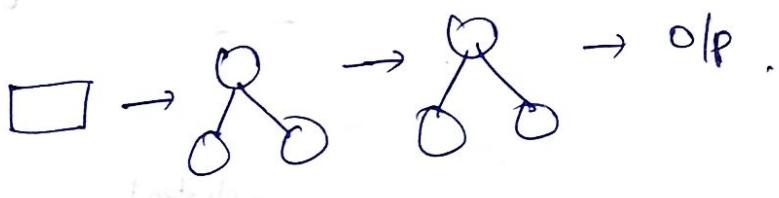
$$\text{Average} = \frac{-11 - 9}{2} = -10$$

$$\alpha_1 = 0.01$$

$SL + \alpha_1 (-10) \Rightarrow O/P$ [But this is for one decision tree]

↓
base model
↓
first decision tree

Like wise Create more decision tree,



$$O/P = SL + \alpha_1(DT_1) + \alpha_2(DT_2) + \alpha_3(DT_3) + \dots + \alpha_n(DT_n)$$

This is about the XGBoost Regressor.

It is also comes under the blackbox model.