

Introduction to Boosting

(Bootstrap Aggregation)

(Base learner) Bagging



→ Random forest classifier
→ Random forest regressor.

Ensembles

{ Bagging: base learner (Decision Trees)

Boosting: Weak learners (Decision trees)

① AdaBoost

② Gradient Boosting

③ XGBoost

=
⇒ for both regression & classification problem statement.

DT → lead to overfitting (high variance)

To reduce this

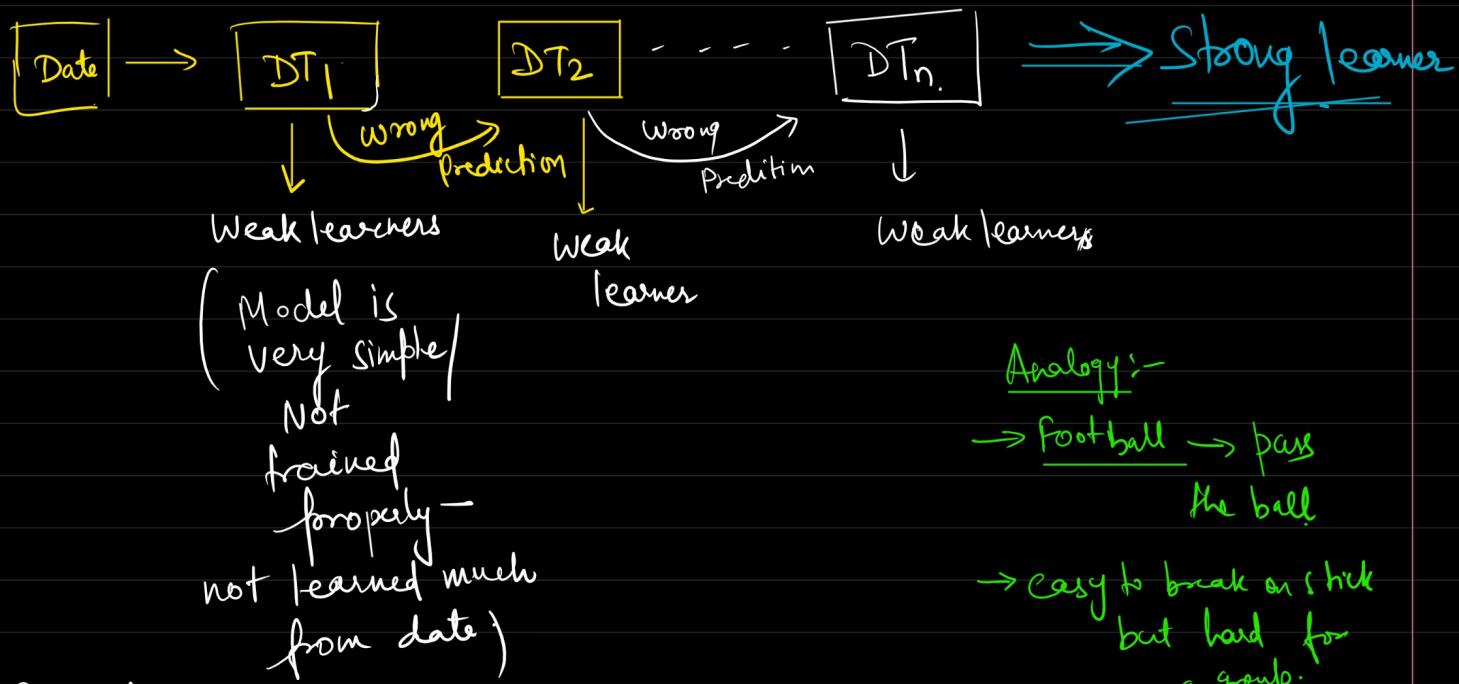
Bagging → Random forest

* Bagging reduces the variance

(multiple overfitted model combining reduces the variance)

* Boosting (Boost \downarrow)

Principles:- Build a first model on training dataset and then build a second model to rectify the errors present in the first model. This procedure is continued until and unless the errors are minimized | data is correctly predicted.



* Prediction

Bagging

- Classification: Voting
- Regression: Average

Boosting

$$\begin{aligned}
 M_1 &\rightarrow M_2 \rightarrow M_3 \rightarrow \dots \rightarrow M_n \\
 \alpha_1(M_1) + \alpha_2(M_2) + \dots + \alpha_n(M_n) & \\
 \downarrow & \\
 \text{Strong learners} &
 \end{aligned}$$

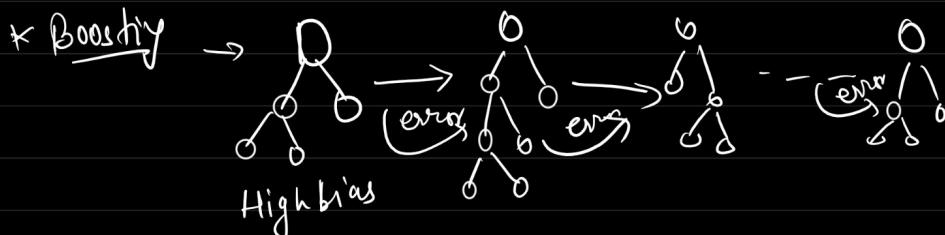
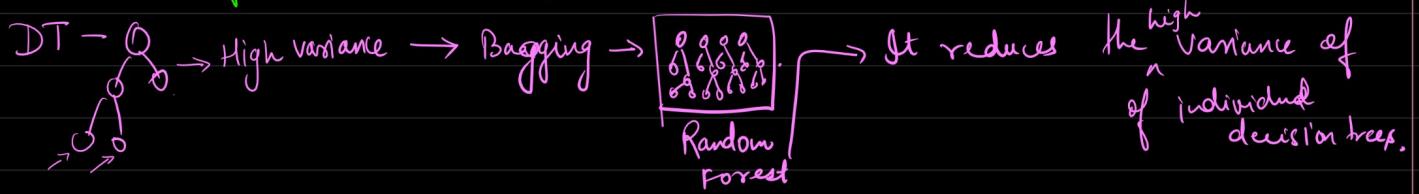
$\alpha_1, \alpha_2, \dots, \alpha_n \Rightarrow \underline{\text{weights}}$

Boosting is a machine learning ensemble technique that improves model accuracy by combining multiple weak learners, typically decision trees. It operates sequentially, where each model corrects the errors of its predecessor by giving more weight to misclassified instances. Popular algorithms include AdaBoost, Gradient Boosting, XGBoost, LightGBM, and CatBoost. Boosting's strength lies in its ability to enhance predictive performance, making it effective for various applications. However, it requires careful hyperparameter tuning to avoid overfitting. Boosting is widely used due to its capacity to create robust models from relatively simple learners.

AdaBoost (Adaptive Boosting)

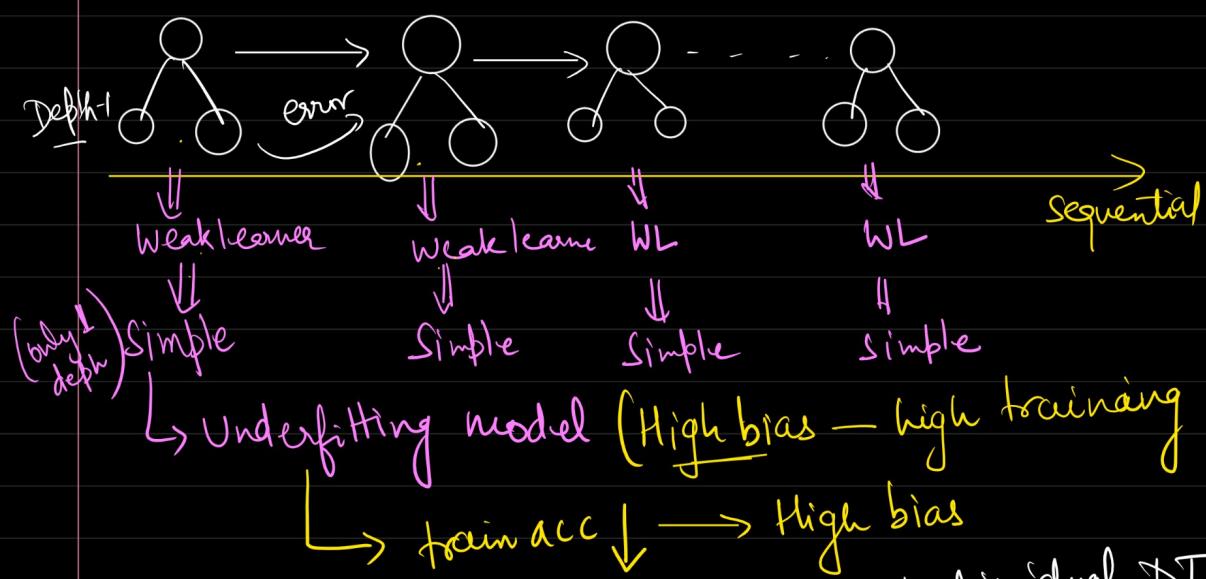
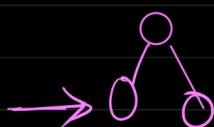
by Yaav Freund and Robert Sopher
in 1995.

- * Works both for classification & regression problem



- * AdaBoost follows the concept of boosting.

Decision Stumps :- A decision tree with only one level of depth



* Each individual DT is a high bias model but error is transferred to next model leading to low bias at an overall level.

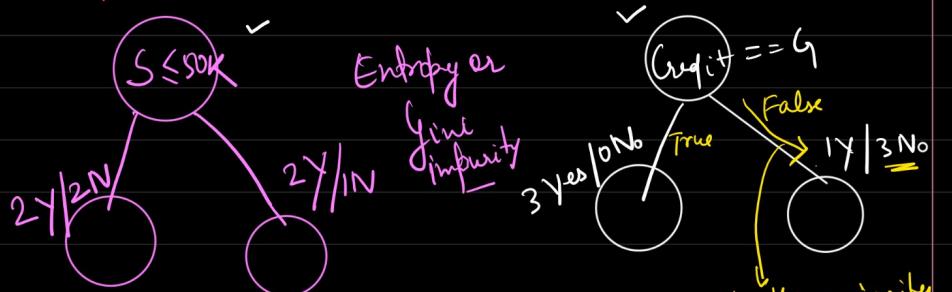
- AdaBoost solves the problem of high bias.
- AdaBoost is low bias model.

→ We always want low bias & low variance model.

* Mathematical explanation of AdaBoost

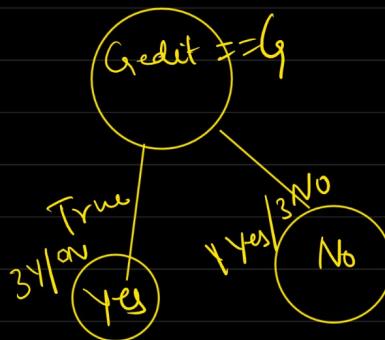
Step-1 → Create decision tree stump.

Salary	Credit Score	Approved.
<= 50K	B	No
<= 50K	G	Yes
<= 50K	G	Yes
> 50K	B	No
> 50K	G	Yes
> 50K	N	Yes
≤ 50K	N	No



Step-2 → Assign equal weights to all of the data points and predict using the higher IG Decision Stump

Salary	Credit Score	Approved.	Sample weights
<= 50K	B	No	$\frac{1}{7}$
<= 50K	G	Yes	$\frac{1}{7}$
<= 50K	G	Yes	$\frac{1}{7}$
> 50K	B	No	$\frac{1}{7}$
> 50K	G	Yes	$\frac{1}{7}$
> 50K	N	Yes	$\frac{1}{7}$
≤ 50K	N	No	$\frac{1}{7}$



Step-3 → Total Error and performance of stump.

$$\text{Total Error} = \frac{1}{7}$$

$$*\text{Performance of Stump} = \frac{1}{2} \ln \left[\frac{1 - TE}{TE} \right]$$

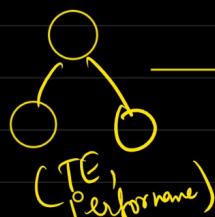
$$f = \alpha_1 M_1 + \alpha_2 M_2 + \dots + \alpha_n M_n$$

$\alpha_1 = 0.896$ → weight of first decision stump.

$$= \frac{1}{2} \ln \left[\frac{1 - \frac{1}{7}}{\frac{1}{7}} \right]$$

$$= \frac{1}{2} \ln (6) \approx 0.896$$

base to e

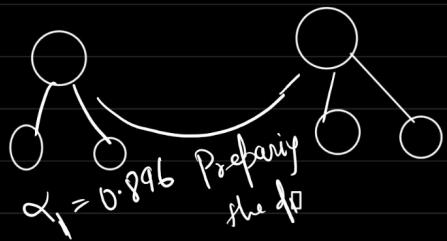


next step is to give more priority to wrong prediction.

Salary	Credit Score	Approval.	Sample weights	Step 4 → Give more weight to incorrect classified datapoint and lesser weight to correct datapoint
$\leq 50k$	B	No	$y_7 \downarrow$	
$\leq 50k$	G	Yes	$y_7 \downarrow$	
$\leq 50k$	G	Yes	$y_7 \downarrow$	
$> 50k$	B	No	$y_7 \downarrow$	* For correct classified dP's - Performance of Stump
$> 50k$	G	Yes	$y_7 \downarrow$	
$> 50k$	N	Yes	$y_7 \uparrow$	$= \text{weight} * e^{-0.896} = 0.058 \quad e=2.73$
$\leq 50k$	N	No	$y_7 \downarrow$	* Incorrect classified dP - + Performance of Stump.

Step-5 Normalized wt & Bin assigned.

Salary	Credit Score	Approval.	SW (init)	Updated wt	Normalized wt	Bin assignment
$\leq 50k$	B	No	$y_7 \downarrow$	0.058	0.08	0-0.08
$\leq 50k$	G	Yes	$y_7 \downarrow$	0.058	0.08	0.08-0.16
$\leq 50k$	G	Yes	$y_7 \downarrow$	0.058	0.08	0.16-0.24
$> 50k$	B	No	$y_7 \downarrow$	0.058	0.08	0.24-0.32
$> 50k$	G	Yes	$y_7 \downarrow$	0.058	0.08	0.32-0.40
$> 50k$	N	Yes	$y_7 \uparrow$	0.349	0.50	0.40-0.90
$\leq 50k$	N	No	$y_7 \downarrow$	0.058	0.08	0.90-0.98



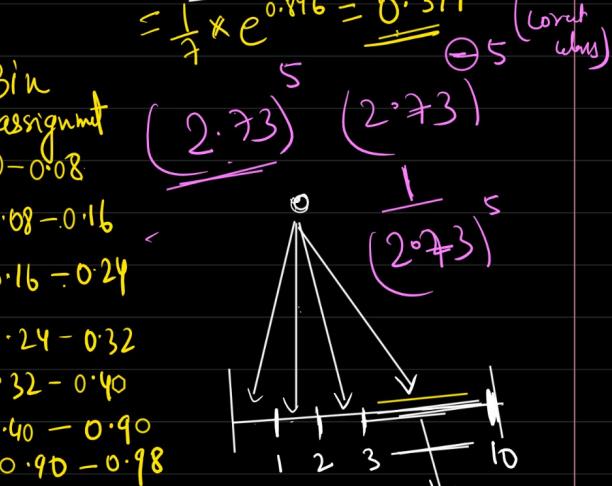
Salary	Credit Score	Approval.	SW	Updated wt	Normalized wt	Bin assignment
$\leq 50k$	B	No	$y_7 \downarrow$	0.058	0.08	0-0.08
$\leq 50k$	G	Yes	$y_7 \downarrow$	0.058	0.08	0.08-0.16
$\leq 50k$	G	Yes	$y_7 \downarrow$	0.058	0.08	0.16-0.24
$> 50k$	B	No	$y_7 \downarrow$	0.058	0.08	0.24-0.32
$> 50k$	G	Yes	$y_7 \downarrow$	0.058	0.08	0.32-0.40
$> 50k$	N	Yes	$y_7 \uparrow$	0.349	0.50	0.40-0.90
$\leq 50k$	N	No	$y_7 \downarrow$	0.058	0.08	0.90-0.98

Step-6 Generate a random no between $0 \text{ to } 1$

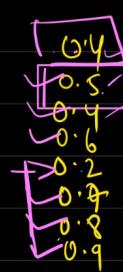
- 0.3
- 0.4
- 0.5
- 0.6
- 0.7
- 0.8
- 0.9

Here total weight is not 1, then we have to normalize

Higher prob of stone falling in this region



Salary	Credit score	Approved	Bin assignment	<u>in bin</u>
$\leq 50k$	B	No	0-0.08	0.4
$\leq 50k$	G	Yes	0.08-0.16	0.5
$\leq 50k$	G	Yes	0.16-0.24	0.6
$> 50k$	B	No	0.24-0.32	0.7
$> 50k$	G	Yes	0.32-0.40	0.8
$> 50k$	N	Yes	0.40-0.90	0.9
$\leq 50k$	N	No	0.90-0.98	0.9



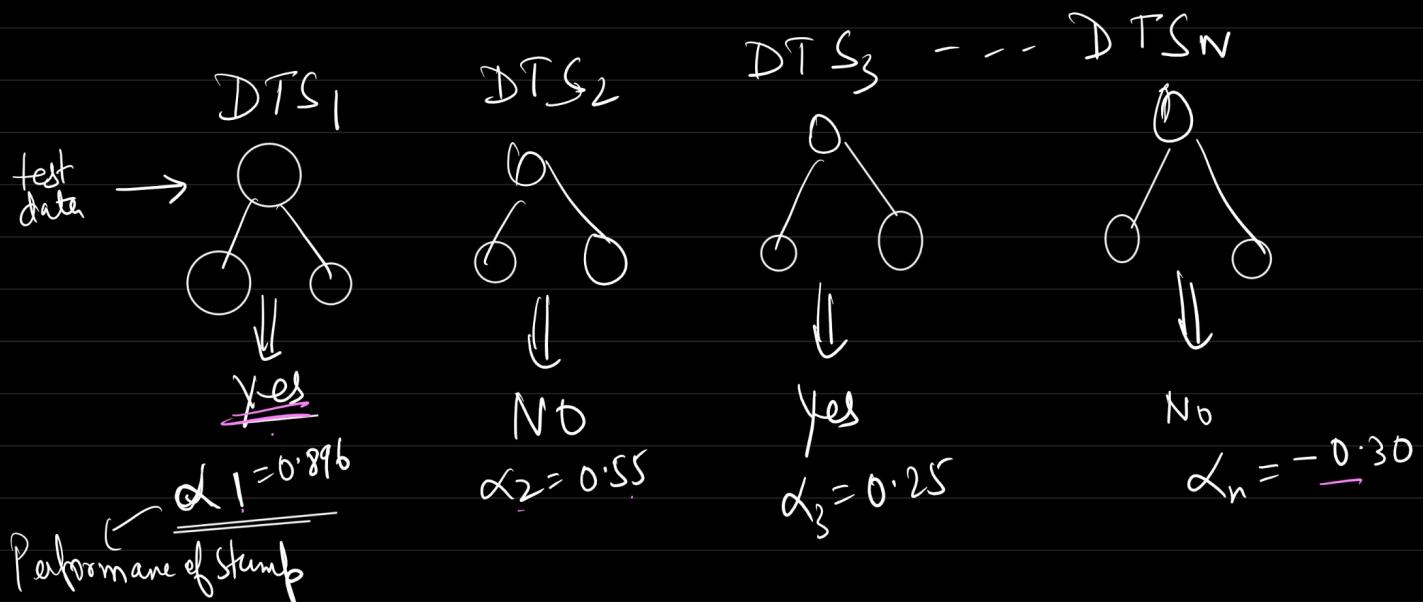
* These records will be sent to next Decision Stump

* Note: → The row with higher bin range due to misclassification can be repeat ⇒ more priority to misclassified df.



* Final Prediction

Test data ($\leq 50k, G$)



$$\begin{aligned}
 f &= \alpha_1(M_1) + \alpha_2(M_2) + \alpha_3(M_3) \dots \alpha_4(M_4) \\
 &= 0.896(\text{Yes}) + 0.55(\text{No}) + 0.25(\text{Yes}) \quad - - - 0.30(\text{No}) \\
 &= 1.15(\text{Yes}) + 0.25(\text{No}) \rightarrow \text{Performance of No} = 0.25 \\
 &\qquad\qquad\qquad \text{final output} \Rightarrow \underline{\text{Yes}}. \\
 &\text{Performance of Yes} = 1.15
 \end{aligned}$$

* multiclass \rightarrow same method \rightarrow there will three Aggregated alpha.

AdaBoost Regressor

→ All the steps will be same except instead of Information gain, you will use Variance reduction to select the decision tree stump.

→ during prediction

$$f = \alpha_1 M_1 + \alpha_2 M_2 + \dots + \alpha_n M_n$$

$M_1, M_2, \dots, M_n \Rightarrow$ will be or continuous Value.

AdaBoost (Adaptive Boosting) is an ensemble learning algorithm that combines multiple weak classifiers to form a strong classifier. It iteratively adjusts the weights of training data points, emphasizing those misclassified by previous classifiers. Each weak classifier is trained to correct errors from the previous ones, and their contributions are weighted based on accuracy. The final model is a weighted vote of these classifiers. AdaBoost is known for improving classification accuracy and reducing overfitting, but it can be sensitive to noisy data and requires more computational resources due to multiple training iterations.

Steps of AdaBoost Algorithm:

Initialize weights: Assign equal weights to all training examples.

Train weak classifier: Train a weak classifier on the weighted training data.

Evaluate classifier: Calculate the classifier's error rate.

Update weights: Increase the weights of misclassified examples and decrease the weights of correctly classified ones.

Calculate classifier weight: Determine the weight of the classifier based on its accuracy.

Combine classifiers: Add the weak classifier to the final model, weighted by its accuracy.

Repeat: Repeat the process for a specified number of iterations or until the model achieves a desired level of accuracy.

Gradient Boosting Algorithms

① Regression

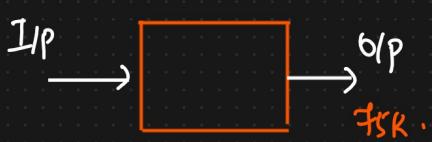
② Classification

Gradient boosting is a machine learning technique used for regression and classification. It builds models iteratively by combining multiple weak learners, typically decision trees, to form a strong predictive model. Each new model is trained to correct the errors of the combined previous models by minimizing a loss function, often using gradient descent. The process involves adding the new model's predictions to the existing model with a controlled learning rate. Gradient boosting is known for its high accuracy and ability to handle a variety of data types but can be computationally intensive and prone to overfitting without proper regularization.

Dataset		y	$(y - \hat{y})$	Predicted \hat{R}_2	\hat{y}	assume dt_1 predicted and then we have calculated residual of it
<u>Exp</u>	<u>Degree</u>	Salary	R_1	R_2	R_3	
→ 2	B.E	50K	-25K	-23K	74.77	-24.77
→ 3	Masters	70K	-5K	-3K	74.97	-4.97
5	Masters	80K	5K	3	-	-
6	PhD	100K	25K	20K	-	-
		75K				

Steps

① Create a Basic Model

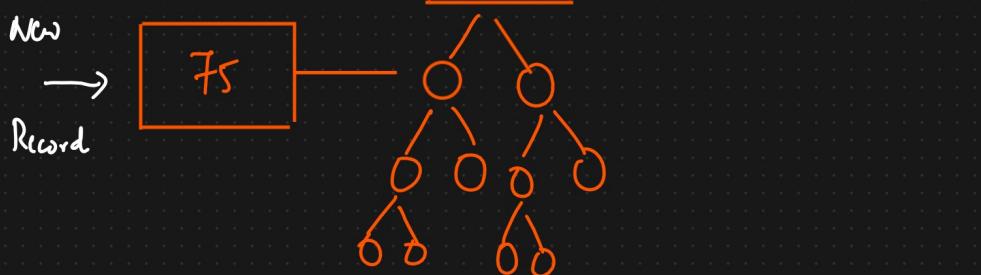


$$\text{Average} = \frac{50K + 70K + 80K + 100K}{4} = 75K$$

② Compute Residuals, Errror

③ Construct a Decision Tree consider inputs x_i and o/p R_i

DT1 $\{x_i, R_i\}$



Basic Model

$$\text{Predicted O/P} = \underset{\downarrow}{75} + \underset{\downarrow}{\text{DT1}} (-23) = 75 - 23 = 52 \{ \text{Overfitting} \}$$

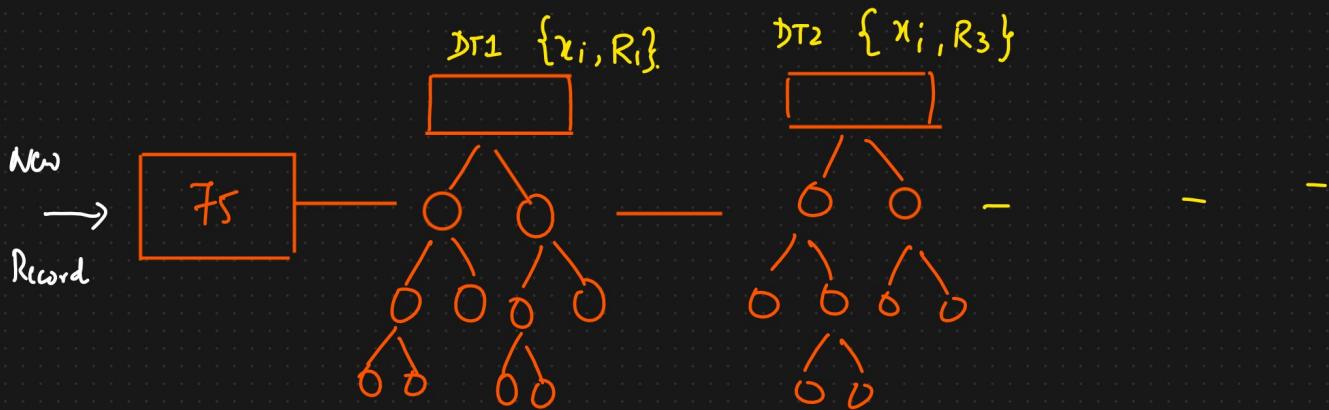
$$\begin{aligned}\text{Predicted O/P} &= \underset{\downarrow}{75} + \underset{\downarrow}{\alpha} (-23) \\ &= 75 + (0.01)(-23) \\ &= 75 - 0.23 \\ &= 74.77\end{aligned}$$

α = Learning Rate

$$\alpha = 0.01 \rightarrow \{0 \text{ to } 1\}$$

Learning Rate: A scaling factor to control the contribution of each weak learner.

$$\begin{aligned}\textcircled{1} \quad \text{Predicted O/P} &= \underset{\downarrow}{75} + \underset{\downarrow}{\alpha} (-3k) \\ &= 75 + (0.01)(-3) \\ &= 75 - 0.03 \\ &= 74.97\end{aligned}$$



Basic Model

$$\text{O/P} = \underset{\downarrow}{75} + \underset{\downarrow}{\alpha_1} (\text{DT}_1) + \underset{\downarrow}{\alpha_2} (\text{DT}_2) + \dots + \underset{\downarrow}{\alpha_n} (\text{DT}_n)$$

Mathematical

$$\alpha_0 = 1$$

$$F(x) = \alpha_0 h_0(x) + \alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots + \alpha_n h_n(x)$$

$$\{\alpha_0, \alpha_1, \alpha_2, \dots, \alpha_n\} \rightarrow \text{Learning Rate } \{0 \text{ to } 1\}$$

$$F(x) = \sum_{i=0}^n \alpha_i h_i(x)$$

Initialization: Start with a base model, usually a simple model like the mean of the target variable for regression tasks.

Iterative Improvement: In each iteration, a new weak learner (typically a decision tree) is trained to predict the residuals (errors) of the combined model from the previous iteration.

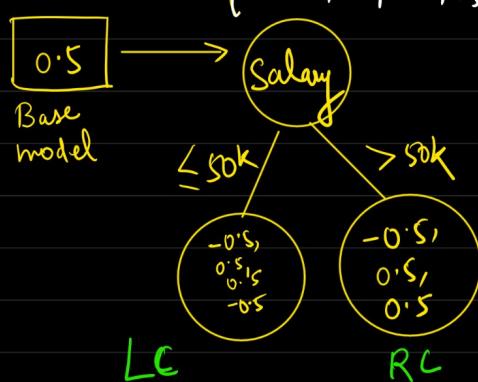
Update Model: The predictions of the new learner are added to the model with a certain weight (learning rate) to update the combined model.

XGBoost (Extreme Gradient Boosting)

→ XGBoost classifier

→ XGBoost Regressor

	Salary	Credit Score	Approval	Base model	R_1
	(y_{act})	(y_{baseP})		($y_{\text{act}} - y_{\text{baseP}}$)	
→ $\leq 50k$	B	0	0.5	-0.5	
$\leq 50k$	G	1	0.5	0.5	
$\leq 50k$	G	1	0.5	0.5	
$> 50k$	B	0	0.5	-0.5	
$> 50k$	G	1	0.5	0.5	
$> 50k$	N	1	0.5	0.5	
$\geq 50k$	N	0	0.5	-0.5	



Step-1 → Create a base model

* Since base model should be unbiased, take 0.5 here

Step-2 - Construct a decision tree with root

Step-3 calculate similarity weights

$$\text{Similarity weight} = \frac{\sum (\text{Residual})^2}{\sum \text{Pr}(\cdot) \text{Pr}(\cdot) + \frac{1}{n} - 1}$$

$$\begin{aligned} * SW_{\text{root}} &= \frac{[-0.5+0.5+0.5-0.5+0.5+0.5-0.5]^2}{0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5)} \\ &= \frac{(0.5)^2}{0.25 + 0.25 + 0.25 + 0.25 + 0.25 + 0.25} \\ &= \frac{0.25}{1.5} = 0.14 \end{aligned}$$

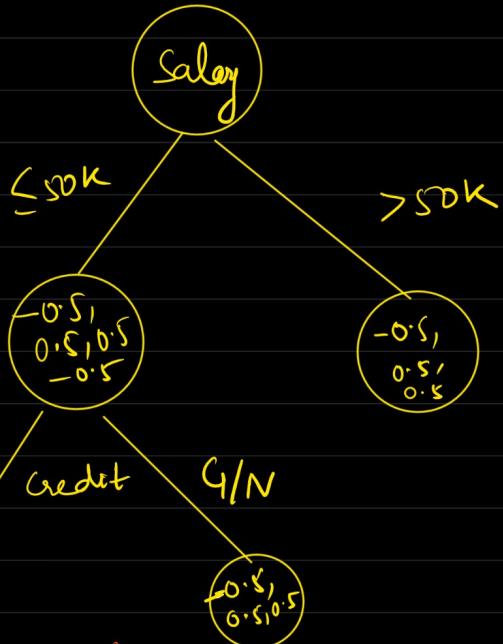
$$* SW_{\text{Left child}} = \frac{(-0.5+0.5+0.5-0.5)^2}{0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5)} = 0$$

$$\begin{aligned} * SW_{\text{Right child}} &= \frac{(-0.5+0.5+0.5)^2}{0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5)} \\ &= \frac{0.25}{0.75} = \frac{1}{3} = 0.33. \end{aligned}$$

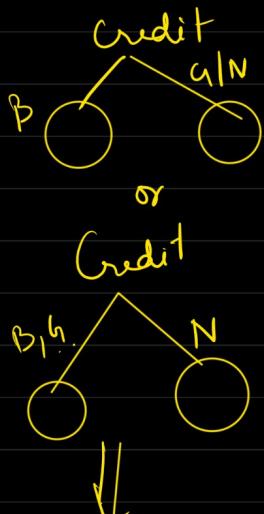
$$\begin{aligned} SW_{\text{gain}} &= (\text{Left child} + \text{Right child}) - \text{Root} \\ &= 0 + 0.33 - 0.14 = 0.19 \end{aligned}$$

* Whichever feature has the highest similarity weight gain, the split will happen on that feature!

* Here assuming Salary column is giving Similarity weight gain



Two possibilities



based on similarity gain score

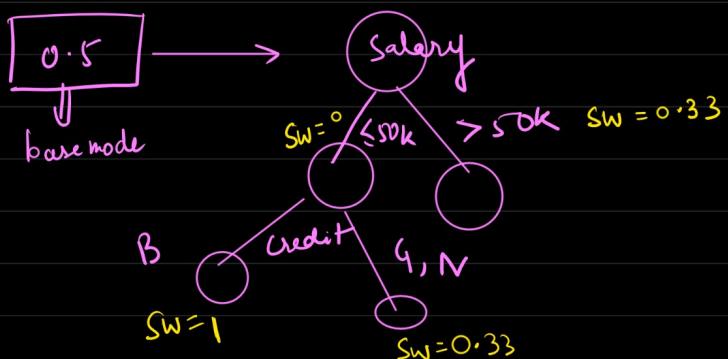
$$SW(LC) = \frac{(-0.5)^2}{0.5(1-0.5)} = 0.25 = 1 \text{ (Pure node)}$$

$$SW(RC) = \frac{-0.5 + 0.5 + 0.5}{0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5)} > 0.25/0.75 = 0.33$$

$$SW(\text{Root}) = \frac{(-0.5 + 0.5 + 0.5 - 0.5)}{(0.25 + 0.25 + 0.25 + 0.25)} = 0$$

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

* How to make prediction ?



XGB classifier

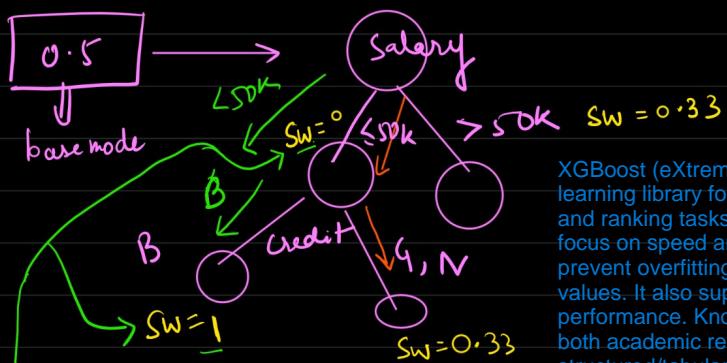
$$O/P = \sigma \left(\text{Base mode} + \alpha(SW) \right)$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

since it's a classification problem, we need to Sigmoid fn.

gives the value b/w 0 & 1

* first row prediction ($\leq 50k, B$)



range of Sigmoid is 0 to 1

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

XGBoost (eXtreme Gradient Boosting) is a high-performance, scalable machine learning library for gradient boosting, which is used for regression, classification, and ranking tasks. It implements decision tree-based ensemble learning with a focus on speed and efficiency. XGBoost includes features like regularization to prevent overfitting, parallel processing, tree pruning, and handling of missing values. It also supports cross-validation and early stopping to improve model performance. Known for its high accuracy and flexibility, XGBoost is widely used in both academic research and industry applications, particularly for structured/tabular data.

$$\begin{aligned} O/P &= \sigma(\text{Base model} + \alpha(\Delta w)) \\ &= \sigma\left(\log\left(\frac{0.5}{1-0.5}\right) + \alpha(0+1)\right) \\ &= \sigma\left(0 + \underset{\text{say } -0.1}{0.1 \times 1}\right) \\ &= \frac{1}{1+e^{-0.1}} = \underline{0.52} \end{aligned}$$

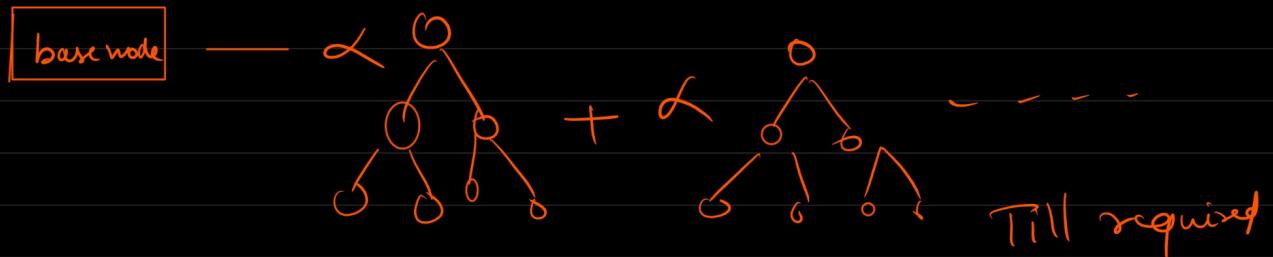
2nd row
($\leq 50k, G$)

$$\begin{aligned} O/P &= \sigma(0 + \alpha(0.33)) \\ &\quad \alpha = 0.1 \\ &= \frac{1}{1+e^{-0.33}} = 0.568 \end{aligned}$$

Step 4 Now do this for all rows calculate \hat{y}

	Credit Score	Approv	Base model ($y_{act} - y_{baseP}$)	\hat{y}
> <= 50k	B	0	0.5 -0.5	0.52
<= 50k	G	1	0.5 0.5	0.508
<= 50k	G	1	0.5 0.5	-
> 50k	B	0	0.5 -0.5	-
> 50k	G	1	0.5 0.5	-
> 50k	N	0	0.5 0.5	-
<= 50k	N	0	0.5 -0.5	-

Step 5 Again construct DT with all and features and \hat{y} as target variable



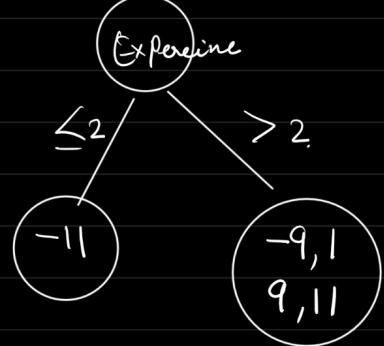
$$O/P = \sigma (\text{Base learner} + \alpha_1(DT_1) + \alpha_2(DT_2) + \dots + \alpha_n(DT_n))$$

of tree is achieved.

* XG Boost Regressor

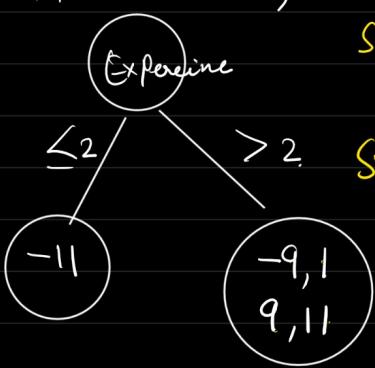
	Exp	Gap	Salary	Base Prediction	<u>Res₁</u>
→ 2	Yes	40	51	-11	
→ 2.5	Yes	42	51	-9	
3	No	52	51	1	
4	No	60	51	9	
4.5	Yes	62	51	11	

(-11, -9, 1, 9, 11)



$$\text{Base Prediction} = \frac{(40+42+52+60+62)}{5} = 51.2 \approx 51$$

(-11, -9, 1, 9, 11)



$$SW_{Root} = \frac{(-11 + -9 + 1 + 9 + 11)^2}{5+1} = \underline{\underline{1}}$$

$$SW_{LC} = \frac{(-11)^2}{1+1} = \frac{121}{2} = 60.5$$

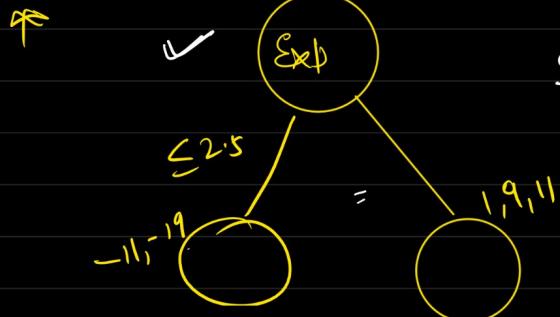
$$SW_{RC} = \frac{(-9 + 1 + 9 + 11)^2}{4+1} \quad \left(\text{Here taking } \underline{\underline{1}} \right)$$

$$= \frac{144}{5} = \underline{\underline{28.8}}$$

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

$$= \underline{\underline{89.134}}$$

* Split will be based on higher SW gain



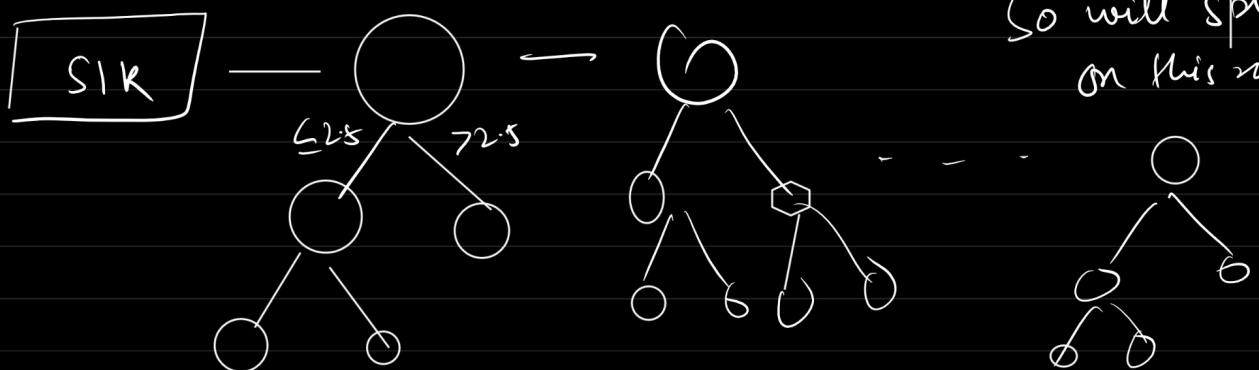
$$SW_{LC} = \frac{(-11 - 9)^2}{2+1} = \frac{400}{3} = 133.3$$

$$SW_{RC} = \frac{(1 + 9 + 11)^2}{3+1} = \frac{21^2}{4} = 110.25$$

$$SW_{Root} = 0.166$$

$$Sw_{gain} = 133.3 + 110.23 - 0.166$$

$\approx 244 \Rightarrow$ Highest gain,
So will split
on this row



until the required
no of dt is made.

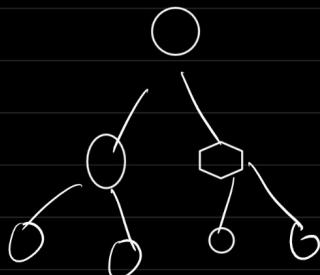
$$\text{Prediction} = 51 + \alpha_1(S) + \alpha_2(\cdot) + \dots + \alpha_n(DT_n)$$

=

$$* \quad \underline{Sw} = \frac{\left(\sum \text{Residue} \right)^2}{\sum p_r(1-p_r) + h} \rightarrow \text{Hyperparameter}$$

One use of h

$$\underbrace{\uparrow \uparrow \underline{Sw} \downarrow \downarrow}_{=}$$



Chances are there

Similarity weight gain can
be negative \rightarrow Stop splitting

$$* \quad \sum p_r(1-p_r) \rightarrow \text{lower value}$$

Any Sw less than $\sum p_r(1-p_r)$ we stop
splitting.

Say $0.5(1-0.5) = 0.25$
if whole Sw is less 0.25, stop splitting.