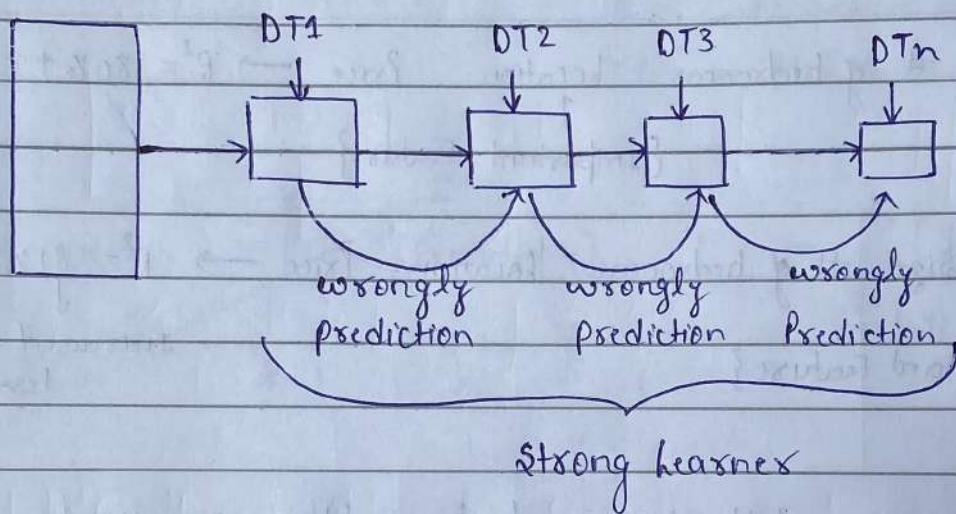


AdaBoost

→ It is a supervised machine learning algorithm which is used to solve both classification and regression problem by using boosting technique.



$$f = \alpha_1(M_1) + \alpha_2(M_2) + \alpha_3(M_3) + \dots + \alpha_n(M_n)$$

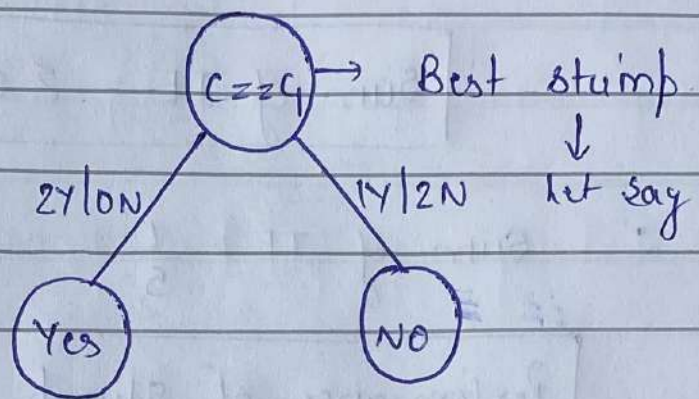
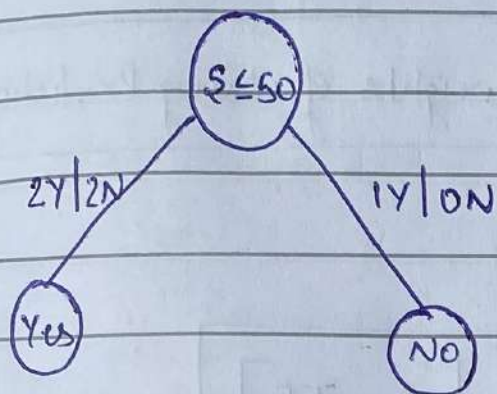
↳ $\alpha \rightarrow$ weights

* Working of AdaBoost :

Dataset

I/P		O/P
Salary	Credit	Approval
$\leq 50k$	B	NO
$\leq 50k$	G	Yes
$\leq 50k$	G	Yes
$> 50k$	N	Yes
$\leq 50k$	N	NO

① Create Decision tree Stump (Decision tree with level 1 only) and select the best stump:



Note: Select best Stump using Entropy or Gini and Information Gain.

② Initialize the sample weights

Salary	credit	Approval	Sample weight	
≤ 50k	B	NO	1/5	
≤ 50k	G	Yes	1/5	
≤ 50k	G	Yes	1/5	
> 50k	N	Yes	1/5	→ wrongly Predicted
≤ 50k	N	NO	1/5	

③ calculate ~~of~~ Sum of total errors and performance of stump.

$$\text{Sum of TE} = \sum \text{sample weights of wrong Prediction}$$

$$\therefore \text{Sum of TE} = \frac{1}{5}$$

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

↓
α.

$$\therefore \alpha_1 = \frac{1}{2} \ln \left[\frac{1 - \frac{1}{5}}{\frac{1}{5}} \right] \approx 0.6$$

④ Update the weights for correctly and incorrectly classified point:

For correct classified point
= weight * e^{α}

$$= \frac{1}{5} \times e^{-0.6} \approx \text{~~0.264~~ } 0.109$$

For incorrectly classified point
= weight * $e^{+\alpha}$

$$= \frac{1}{5} \times e^{+0.6} \approx \text{~~0.264~~ } 0.364$$

Salary	credit	Approval	sw	updated weight
≤ 50k	B	No	1/5	0.109 ↓
≤ 50k	G	Yes	1/5	0.109 ↓
≤ 50k	G	Yes	1/5	0.109 ↓
> 50k	N	Yes	1/5	0.364 ↑
≤ 50k	N	No	1/5	0.109 ↓
				$\Sigma = 0.8$

⑤ Normalize weights computation and assigning bins :

$$\text{Normalize} = \frac{\text{updated weight}}{\sum \text{updated weights}}$$

Normalized Weight

Bins Assignment

0.136

0 - 0.2

0.136

0.2 - 0.3

0.136

0.3 - 0.4

0.455

0.4 - 0.8

0.136

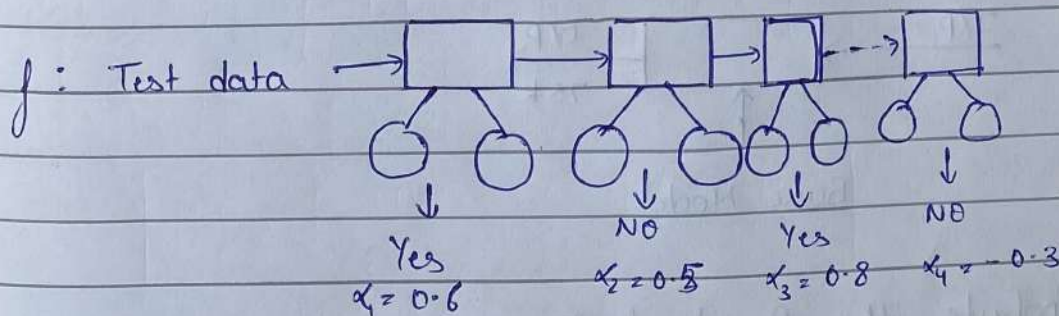
0.8 - 0.9

Biggest Bin
Size
(wrongly predicted
data point)

Note: Now, algorithm will run and iteration process in which it select random numbers between 0 and 1. Most of the numbers will lie inside biggest bin. Hence, the wrongly predicted data point will be selected and send to next decision tree stump.

→ Prediction :

Test data ($\leq 50k, 6$)



$$f = 0.6(\text{Yes}) + 0.5(\text{NO}) + 0.8(\text{Yes}) + (-0.3)(\text{NO})$$

$$= 1.4(\text{Yes}) + 0.2(\text{NO})$$

Performance of say(Yes) ≥ 1.4
Performance of say(No) $= 0.2$
↓
Yes