

# Ensemble Techniques :

## Bagging and Boosting

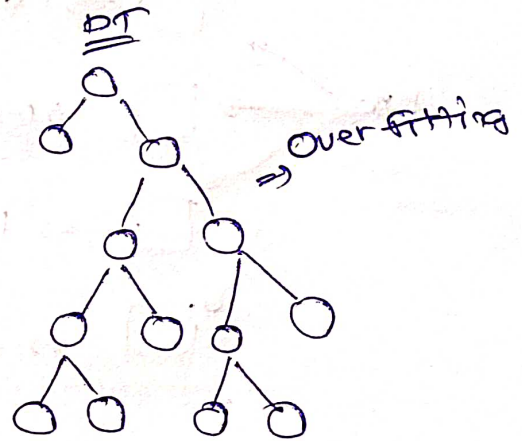
- ① Random Forest Classifier
- ② Random Forest Regression

### Decision Tree :

\* Problem is overfitting

low bias

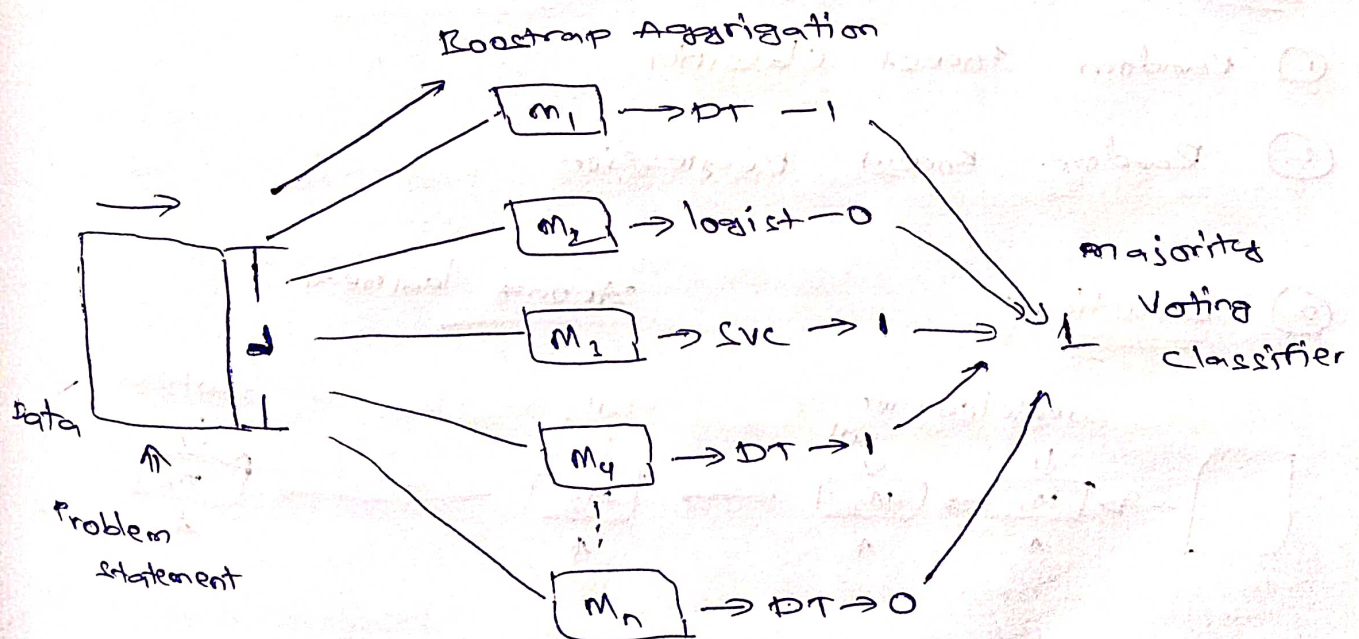
high variance



## Bagging and Boosting

### ① Bagging :

\* Bagging is used when our objective is to reduce the variance of decision tree

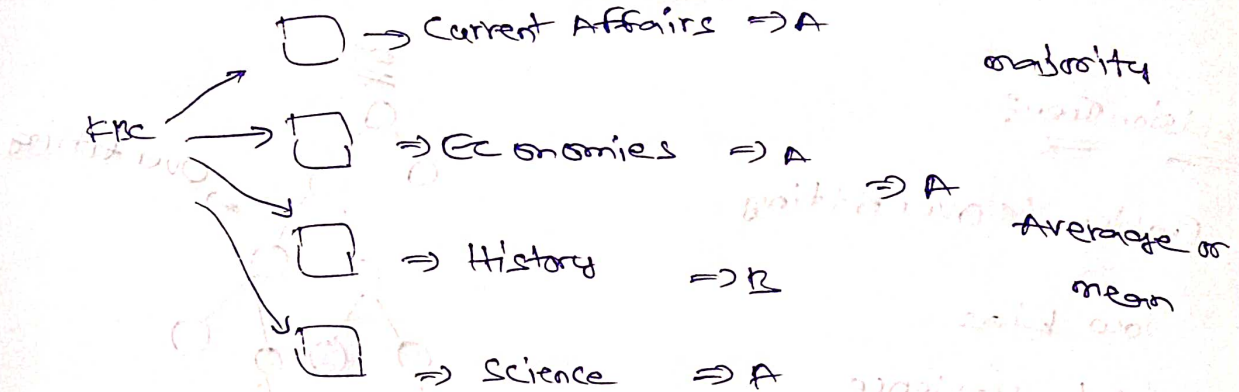


Example:

KBC : Kaun Banega Crorepati (Amitha Bachan)

Krish → Data Science

UPSC → Diversified Apti logical math Soci curr Hist



Used Ensemble Techniques:



Kaggle, HackerRank



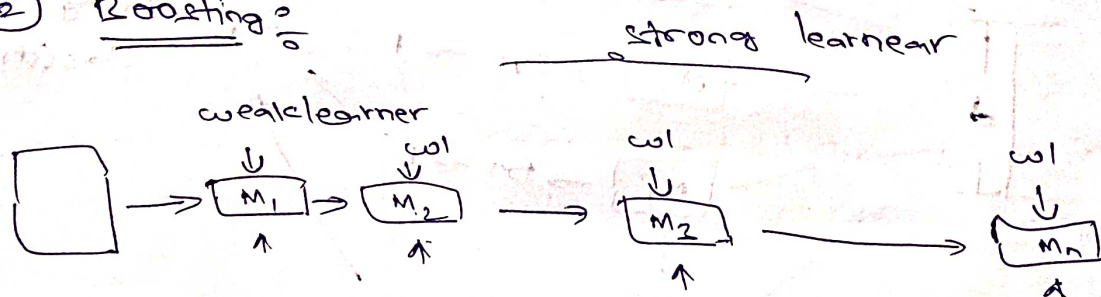
Outperforming results

Bagging

① Random Forest Classifier

② Random Forest Regression

② Boosting



① AdaBoost Regressor And Classifier

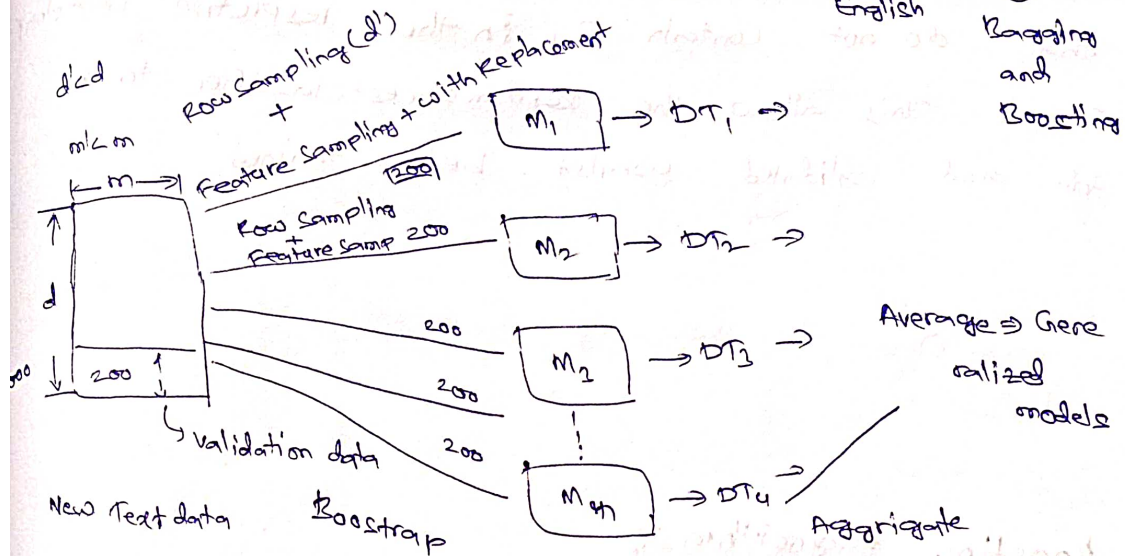
② Gradient Boost Regression and Classifier

② Xgboost Regression and Classifier

↳ xtreem Gradient boost



# Random Forest Classifier and Regression : Ensemble Technique

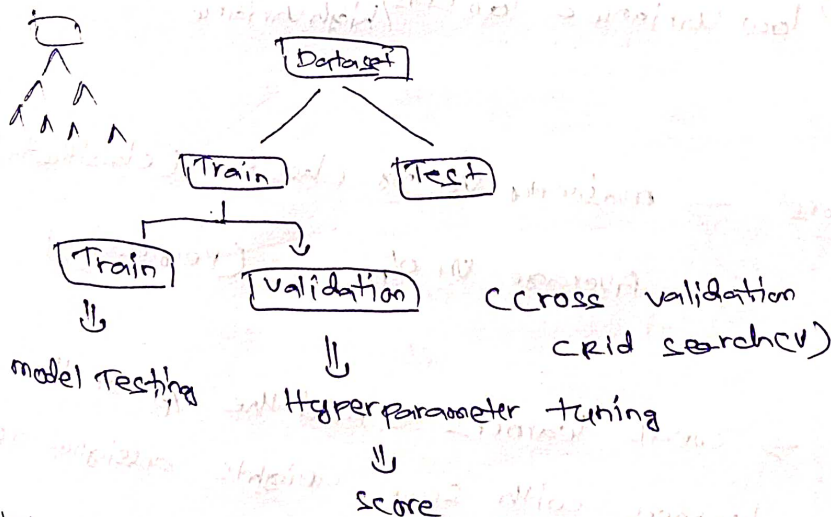


Oob - score = True

DT1  $\rightarrow$  **Overfitting**  $\rightarrow$  **Preprocessing or Post Pre pruning**

① Low Bias  $\rightarrow$  Training

② High Variance  $\rightarrow$  **Low variance**



Oob

out of bag Evaluation

Oob - score  $\Rightarrow$  83%

$\downarrow$   
validation

out of Bag Error

$1 - \text{Oob-score} = 1 - 83 = 0.17\%$

$\downarrow$   
**{validation Error}**

## oob Errors

\* The out of bag (oob) error is the average error for each  $z_i$  calculated using predictions for the trees that do not contain  $z_i$  in their respective bootstrap sample. This allows the Random Forest classifier to be fit and validated whilst being trained.

## Boosting Algorithms:

- \* sequential weak learners
- \* Boosting tries to reduce bias.

### ① Adaboost

step

under-fitting

Boosting  $\left\{ \begin{array}{l} \text{low Bias} \Leftarrow \text{high Bias} \\ \text{low Variance} \Leftarrow \text{low Bias / high Variance} \end{array} \right.$

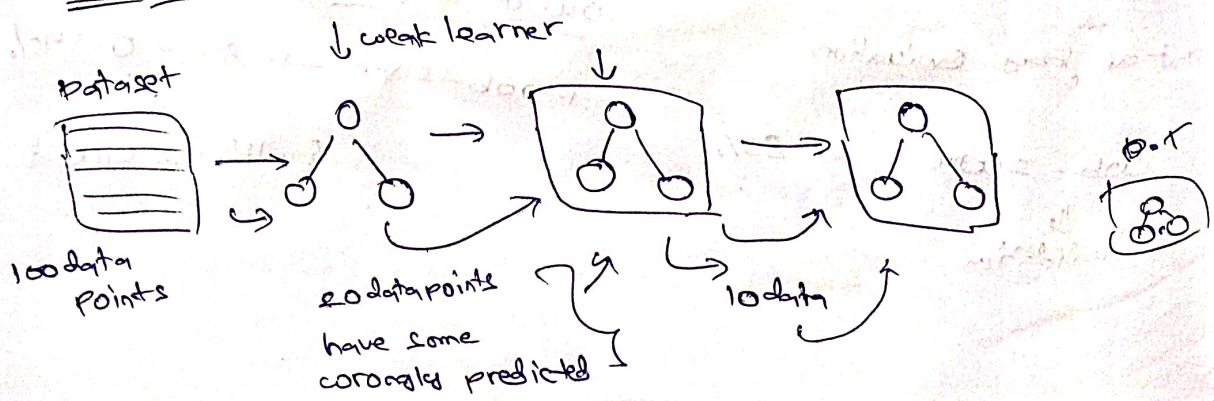
$\left\{ \begin{array}{l} \text{Training data Acc} \uparrow \uparrow \\ \text{Test data Acc} \uparrow \uparrow \text{ 45\%} \end{array} \right.$

Random Forest  $\rightarrow$  majority voting classifier (classification)

Average O/p (Regression)

Adaboost  $\rightarrow$  weak learners  $\rightarrow$  Add the o/p of the weak learners with some weights assigned to it

### Adaboost :





$$f = d_1(M_1) + d_2(M_2) + d_3(M_3) + \dots + d_n(M_n)$$

$m_1, m_2, m_3, \dots, m_n \rightarrow$  weak learners

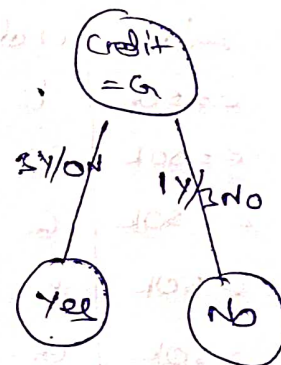
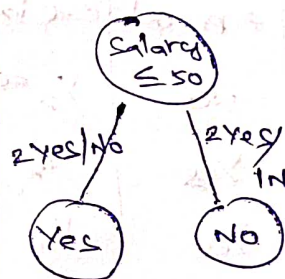
$d_1, d_2, d_3, \dots, d_n \rightarrow$  weights



$\Rightarrow$  total  $\odot$

Salary	Credit	Approval	weights
$\leq 50K$	B	No	$1/7$
$\leq 50K$	G	Yes	$1/7$
$\leq 50K$	G	Yes	$1/7$
$> 50K$	B	No	$1/7$
$> 50K$	G	Yes	$1/7$
$> 50K$	N	Yes	$1/7$
$\leq 50K$	N	No	$1/7$

① We created DT stump by selecting the best one



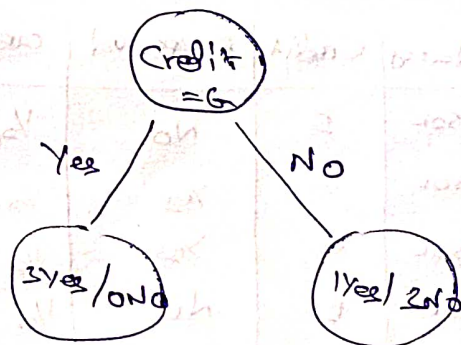
Entropy

$$H(x) = -p + \log_2 p +$$

or Gini

$$-p - \log_2 p -$$

Salary	Credit	Approval	weight	
$\leq 50K$	B	No	$1/7$	$\left\{ \begin{array}{l} B = \text{Bad} \\ G = \text{Good} \\ N = \text{Normal} \end{array} \right.$
$\leq 50K$	G	Yes	$1/7$	
$\leq 50K$	G	Yes	$1/7$	
$> 50K$	B	No	$1/7$	
$> 50K$	G	Yes	$1/7$	
$> 50K$	N	Yes	$1/7$	
$\leq 50K$	N	No	$1/7$	



② Calculate the Total Error (Add the weights of wrong data points)

$$T.E = \frac{1}{7}$$

③ Performance of stump =  $\frac{1}{2} \ln \left( \frac{1 - TE}{TE} \right) = \frac{1}{2} \ln(6) \approx 0.896$

$f = d_1(m_1) + d_2(m_2) + d_3(m_3) + \dots + d_n(m_n)$

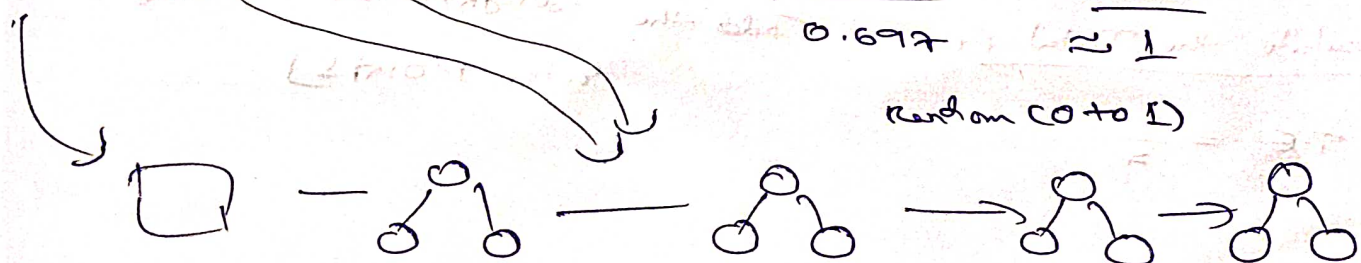
$d_1 = 0.896$

④ Update the weight for Correctly and Incorrectly data points

Salary	Credit	Approved	weight	update weight	
<=50k	B	No	$\frac{1}{2} = 0.14$	0.58	For Correctly Classified points = weight * e <sup>-Performance</sup> = $\frac{1}{2} * e^{-(0.896)}$ = 0.349
<=50k	G	Yes	$\frac{1}{2}$	0.58	
<=50k	G	Yes	$\frac{1}{2}$	0.58	
>50k	B	No	$\frac{1}{2}$	0.58	
>50k	G	Yes	$\frac{1}{2}$	0.58	
>50k	N	Yes	$\frac{1}{2}$	0.349	For Incorrect Classifier point ⇒ weight * e <sup>Performance</sup> ⇒ $\frac{1}{2} * e^{(0.896)}$ ⇒ 0.349
<=50k	N	No	$\frac{1}{2}$	0.58	

⑤ Normalized weights and Assigned Bins

Salary	Credit	Approved	weights	Updated weights	Normalized weights	Bins Assigned
<=50k	B	No	$\frac{1}{2} = 0.14$	0.058 / 0.692	0.08	0 - 0.08
<=50k	G	Yes	$\frac{1}{2}$	0.058 / 0.692	0.08	0.08 - 0.16
<=50k	G	Yes	$\frac{1}{2}$	0.058 / 0.692	0.08	0.16 - 0.24
>50k	B	No	$\frac{1}{2}$	0.058 / 0.692	0.08	0.24 - 0.32
>50k	G	Yes	$\frac{1}{2}$	0.058	0.08	0.32 - 0.40
>50k	N	Yes	$\frac{1}{2}$	0.249	0.50	0.40 - 0.90
<=50k	N	No	$\frac{1}{2}$	0.058	0.08	0.90 - 1
				0.692	1	

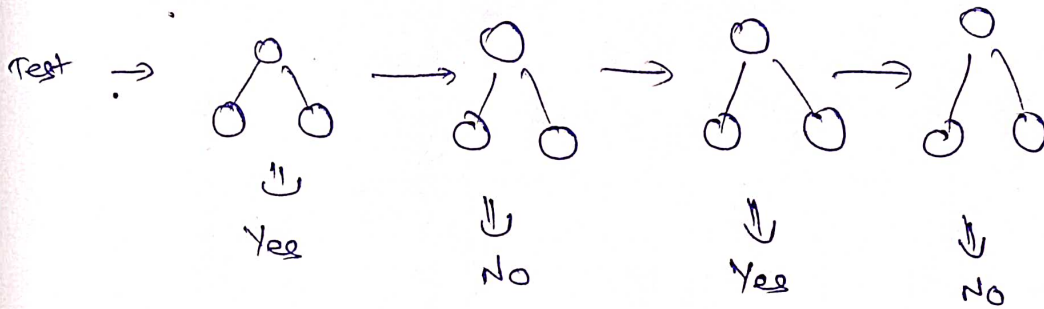




$$d_1(m_1) + d_2(m_2) + d_3(m_3) + \dots + d_n(m_n)$$

### \* Final Prediction

Test ( $\leq 50K$ , G)



$\hookrightarrow$   $d_1 = 0.896$      $d_2 = 0.650$      $d_3 = 0.36$      $d_4 = 0.20$

$$f = L_1(m_1) + L_2(m_2) + L_3(m_3) + L_4(m_4)$$

$$= 0.896(\text{Yes}) + 0.650(\text{No}) + 0.36(\text{Yes}) + 0.20(\text{No})$$

$$= 1.2(\text{Yes}) + 0.85(\text{No})$$

↓  
Yes