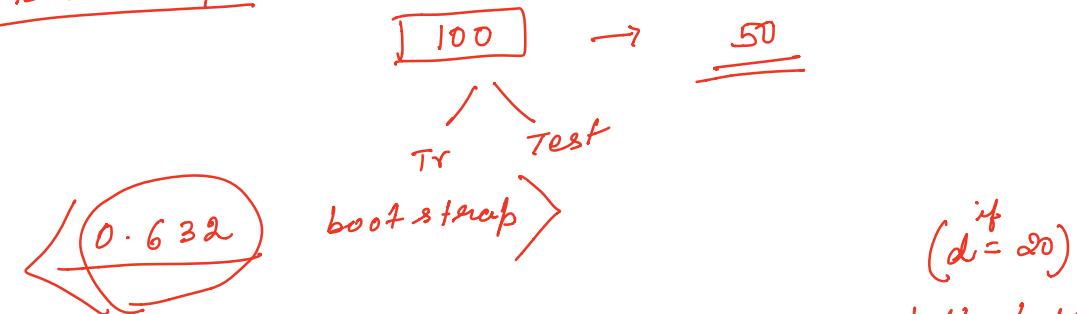


Dec-27

16/03/21

→ Bootstrap  
→ ROC curve.

### Bootstrap.



Suppose we are given a dataset of ' $d$ ' tuples.  
The dataset is sampled ' $d$ '-times with replacement resulting in a bootstrap sample or training set of ' $d$ ' samples.

Trying out several times, it is observed that  $63.2\%$  of original data will end up in bootstrap and remaining  $36.8\%$  will form the test set.

prob. that each tuple will be selected =  $1/d$   
" " " " not being " =  $(1 - 1/d)$

As ' $d$ ' times the process is repeated,

the prob. that tuple not being selected =  $(1 - 1/d)^d$

If  $d$  is large, the prob. approaches  $e^{-1} = 0.368$

Thus,  $36.8\%$  of tuples will not be selected for training and thereby end up in test set. The remaining  $63.2\%$  will form the training set.

→ The sampling procedure is repeated  $K$ -times.

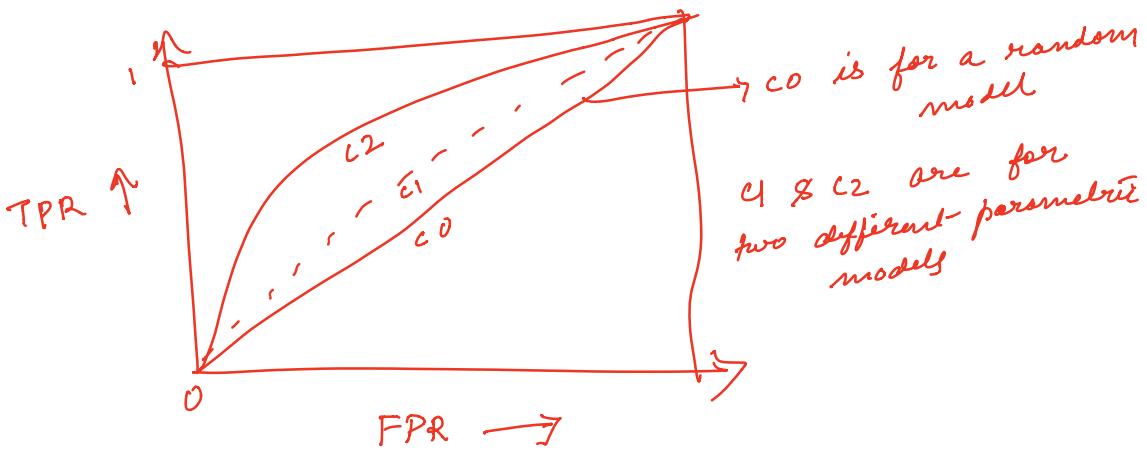
$$Acc(M) = \sum_{i=1}^K (0.632) \times Acc(M_i)_{\substack{\text{test} \\ \text{set}}} + 0.368 \times Acc(M_i)_{\substack{\text{train} \\ \text{set}}}$$

Model accuracy over  
test set

train set

ROC curve : (Receiver operating characteristic curve)

→ The relative increase of FPR with an increase in TPR is captured in ROC of a model.



$(0,0)$  represents all cases to be negative.

$(1,1)$  " " " " " positive

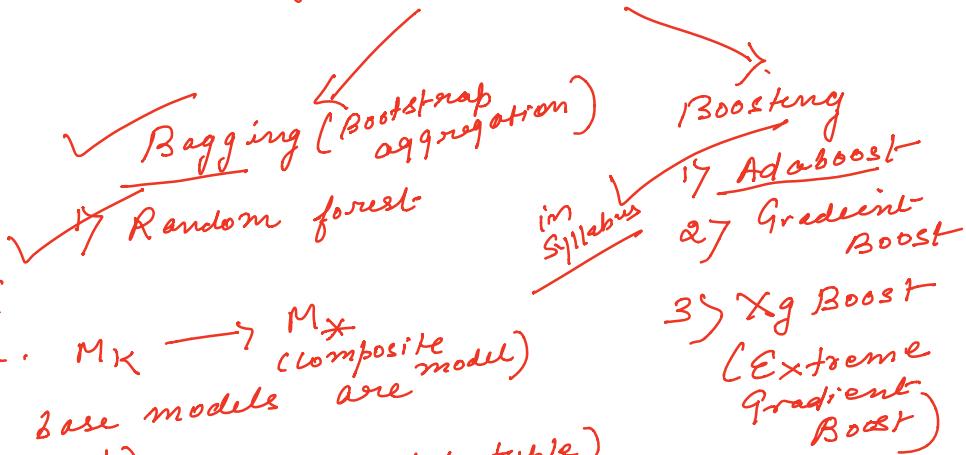
$(0,1)$  " perfect classifier

$(1,0)$  point represents incorrect classifications for every positive that exists.

The diagonal line from  $(0,0)$  to  $(1,1)$  represents a random model where we can expect an equal no. of true & false cases.

8.6  
Ham Karim  
3rd Edn

## Techniques to improve classification accuracy (Ensemble approaches)

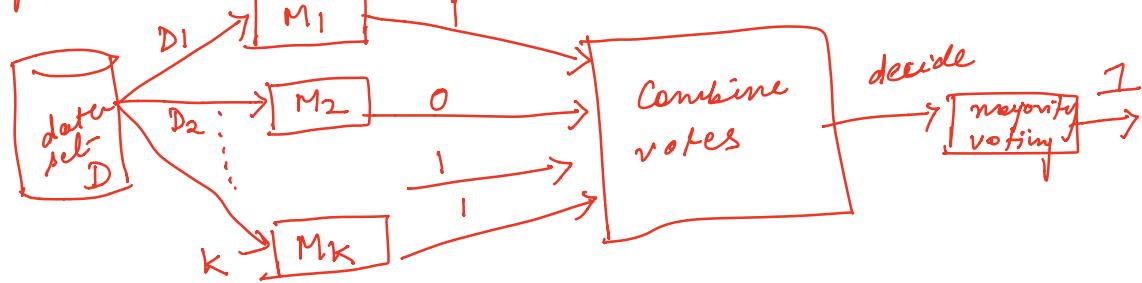


Bagging :

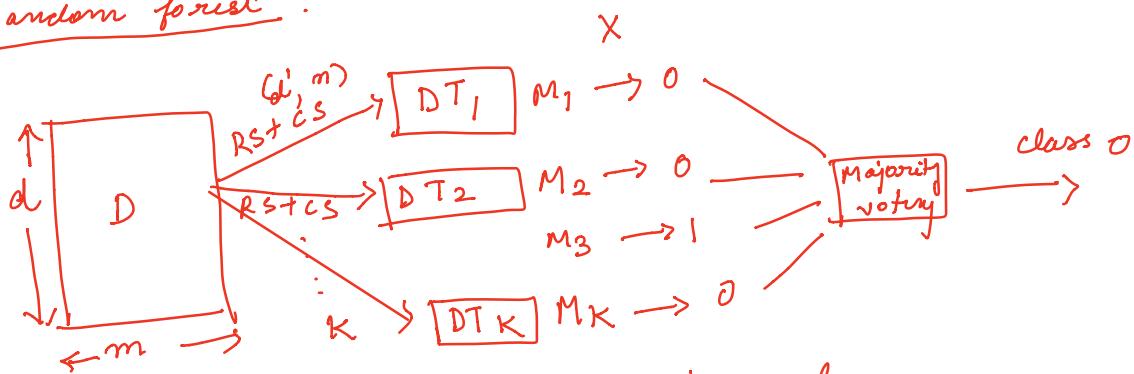
$M_1, M_2, \dots, M_K \rightarrow M^*$  (composite model)

( $K$  no. of base models are created)

(row sampling)

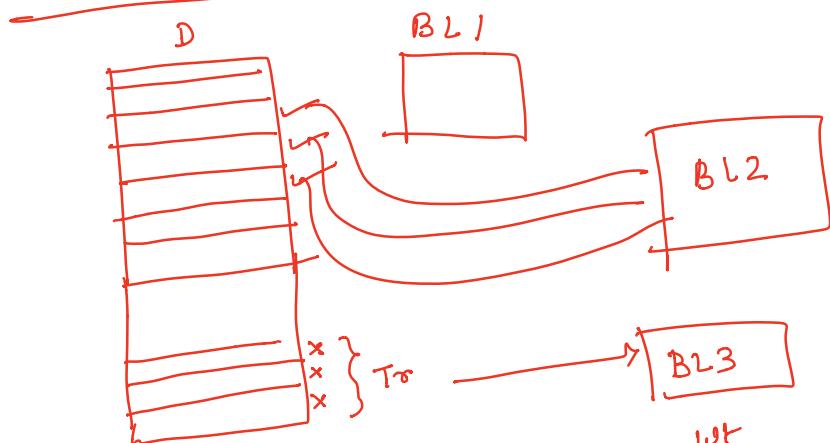


Random forest :



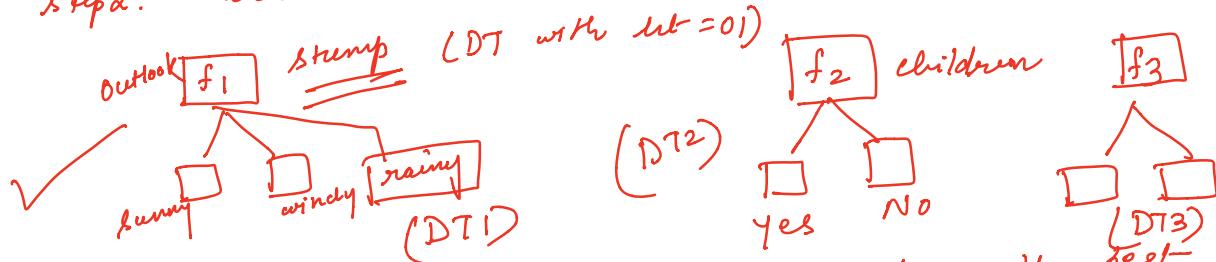
Row sampling ( $d'$  rows) s.t.  $d' \ll d$   
 column sampling ( $n$  columns) s.t.  $n \ll m$

Adaboost : base learner



	$f_1$	$f_2$	$f_3$	Op	Step 1 Sample $w_t = w_{i/n}$	New Wt	Normalize wts
1	✓			0	1/7	0.05	0.05/0.649 → 0.077
2	✓			1	1/7	0.05	0.05/0.649 → 0.077
3	X			0	1/7	0.349	0.349/0.649 → 0.53
4	✓			1	1/7	0.05	:
5	✓			0	1/7	0.05	:
6	✓			1	1/7	0.05	:
7	✓			0	1/7	0.05	0.07
					$\sum = 1$	$\sum = 1$	

Step 2: create the base learners (here it is DT)



(use min entropy, Give to choose the best stump)

Step 3: Evaluate error and performance of stump

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

$$\text{Performance of Stump} = \frac{1}{2} \log_e \left( \frac{1 - \text{TE}}{\text{TE}} \right)$$

$$\begin{aligned}
 &= \frac{1}{2} \log_e \left( \frac{1-1/7}{1/7} \right) \\
 &= \frac{1}{2} \log_e (6) \\
 &= 0.896
 \end{aligned}$$

Step 4 : update wts. s.t increase wts. of misclassified record and decrease wts. of correct classified record.

i) incorrectly classified record

$$\begin{aligned}
 \text{New wt} &= \text{old wt} \times e^{\text{total error}} \\
 &= \frac{1}{7} \times e^{0.895} \\
 &= 0.349
 \end{aligned}$$

ii) correctly classified record

$$\begin{aligned}
 &= \frac{1}{7} \times e^{-(0.895)} \\
 &= 0.05
 \end{aligned}$$

Step 5 : Normalize the wts. so that the summation = 1 ( $\leq \text{wt} = 01$ )

<u>Step 6</u>	<u>create buckets</u>	<u>Generate random no 8 times</u>
	0 - 0.07	0.43 rec 03
	0.07 - 0.14	0.31 rec 03
	0.14 - 0.67 ✓	0.68 rec 04
	0.67 - 0.74	0.58 rec 03
	0.74 - 0.81	
	0.81 - 0.88	
	0.88 - 0.95	

Accordingly, the records are selected and the same process continues till the classification accuracy is attained as per user specified level.