

Module - 3

Ensemble learning →

3.1 Understand Ensembles

K-fold Cross Validation

⇒ Boosting

Stumping (Adaboost)

XGBoost

3.2 ⇒ Bagging

Subagging

Random Forest

Comparison with Boosting

Different ways to

Combine Classifiers

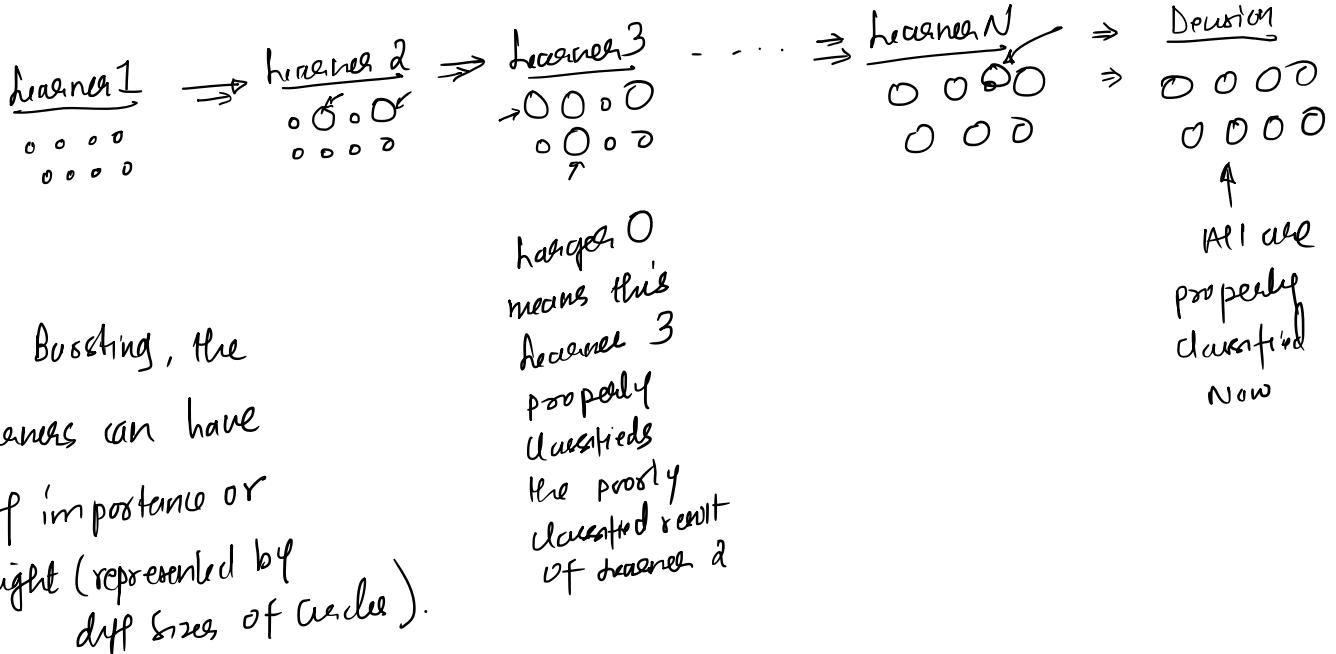
* Ensemble ->

- * In ML, ensemble is a model that combines the prediction from two or more models.
- * The models that contributes to Ensemble are known as ensemble members.
- * The members may or may not be trained on same training data and they may be of same type or different type.
- * It's very powerful method to improve the performance of the model.
- * It's technique that uses group of weak learners in order to create a strong and aggregated learner.

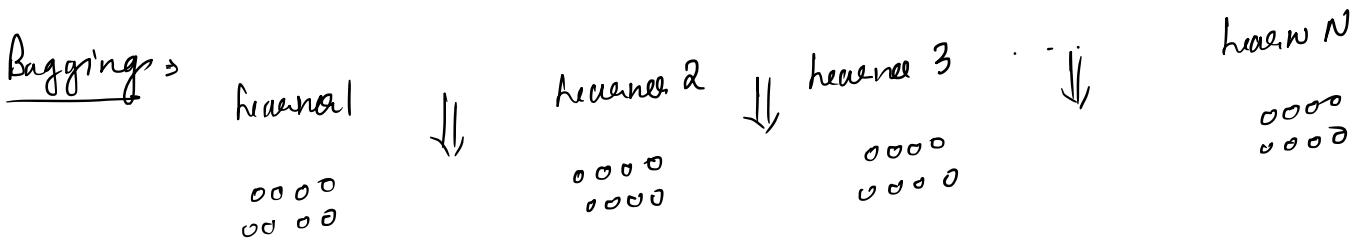
→ The Ensemble technique helps to reduce the Variance (By Bagging) and Bias (By Boosting) and thus helps in improving the predictions.

→ Boosting model :-
* It falls inside family of Ensemble method.
* It consists of filtering or weighting the data that is used to train team of Weak Learners, so that the new learner can give more weight on sample that is poorly classified by previous learner.
→ In Boosting the learners are trained

Sequentially



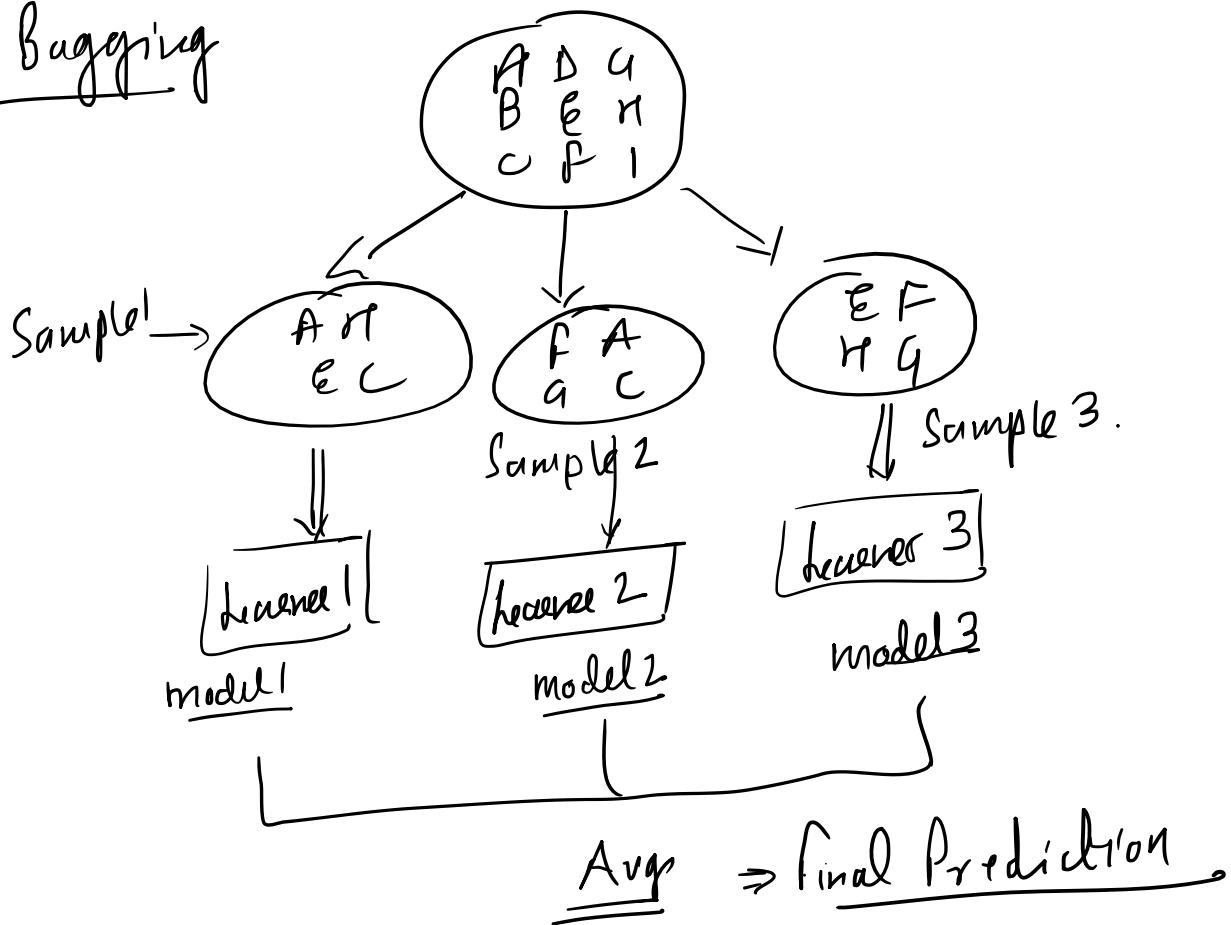
* In Boosting, the learners can have diff importance or weight (represented by diff sizes of circles).



- ⇒ In Bagging the weak learners are trained in parallel using randomness
- ⇒ All learners have same weights.

Note

Bagging

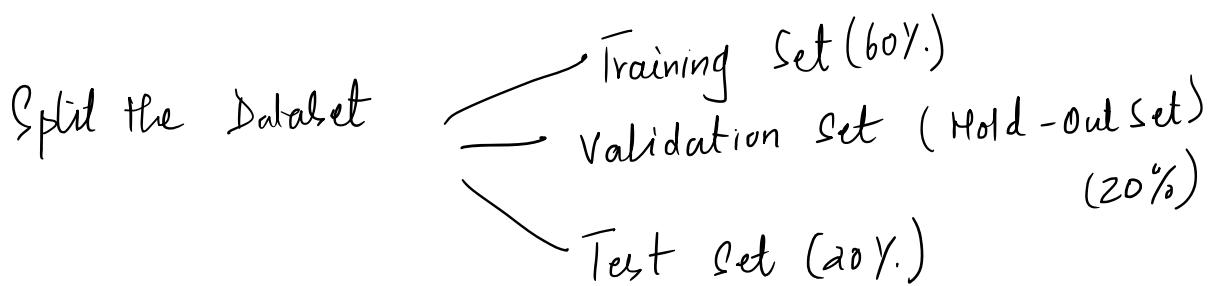


- * Bagging helps in reducing variance
- * Boosting helps in reducing bias.

Cross Validation →

- In Supervised ML
- Train a Model on a Dataset
- Trained model is used to predict the target given new sample.
- How to know if model we have trained will produce effective and accurate result on new input

Cross Validation → It is process that ensures the model will perform well on new Data.



Training Set = part of data on which model is trained.
(This dataset will help to * build model)

* Validation Set ⇒ * Evaluate the Model

- will help to chk if model overfits or Underfits
- update the parameters and again train the model
- Repeat this until the model performs best on Validation set

→ ~~Repeat until convergence~~
Validation set

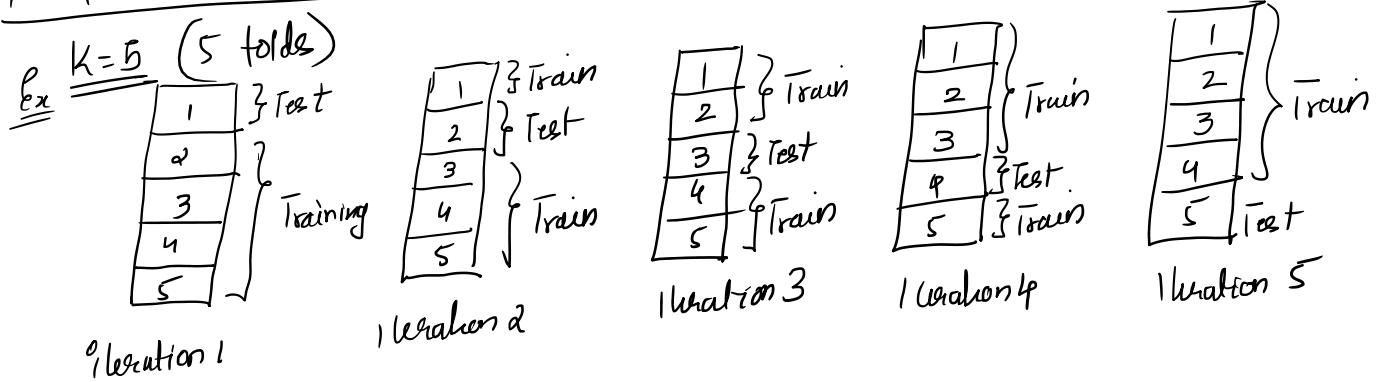
Test Set \Rightarrow^* Prediction

→ The fully trained model after being evaluated on validation set can be used on test set to generate Estimation.

Q) Types of Cross Validation →

- ① The standard validation set Approach
- ② (Leave one out) leave one out cross validation.
- ③ K-fold Cross Validation

K-fold Cross Validation →

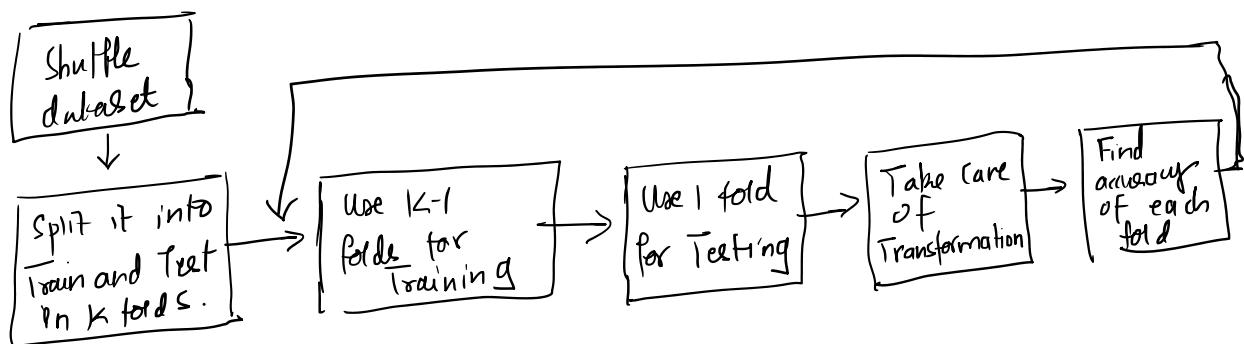


K fold → K fold helps us to build the model in generalized form.

→ To achieve this K fold Cross Validation splits the dataset into Training Testing & Validation.

→ Here Test and Train data will support building the model.

→ Life cycle of K-fold Cross Validation.



- * The No of iterations ideally is K time
- * Finding mean of accuracy score of each iteration will give the consistency of the Trained model.

Rules

$$\textcircled{1} \quad \underline{\underline{k \geq 2}}$$

if $k=2 \rightarrow$ just 2 iteration.

if $k=n \geq 2 \rightarrow$
 $n-1$ for Training
 1 for Testing

\textcircled{2} most commonly used value of $\underline{\underline{k=10}}$

\textcircled{3} If k is very large then the running time of process will increase.

\textcircled{4} The value of k is inversely proportional to size of data i.e if dataset size is small then number of folds can increase.

Bagging \rightarrow Random Forest

* Random Forest is example of Bagging.

* You must know Decision Tree Construction [Gini Index]

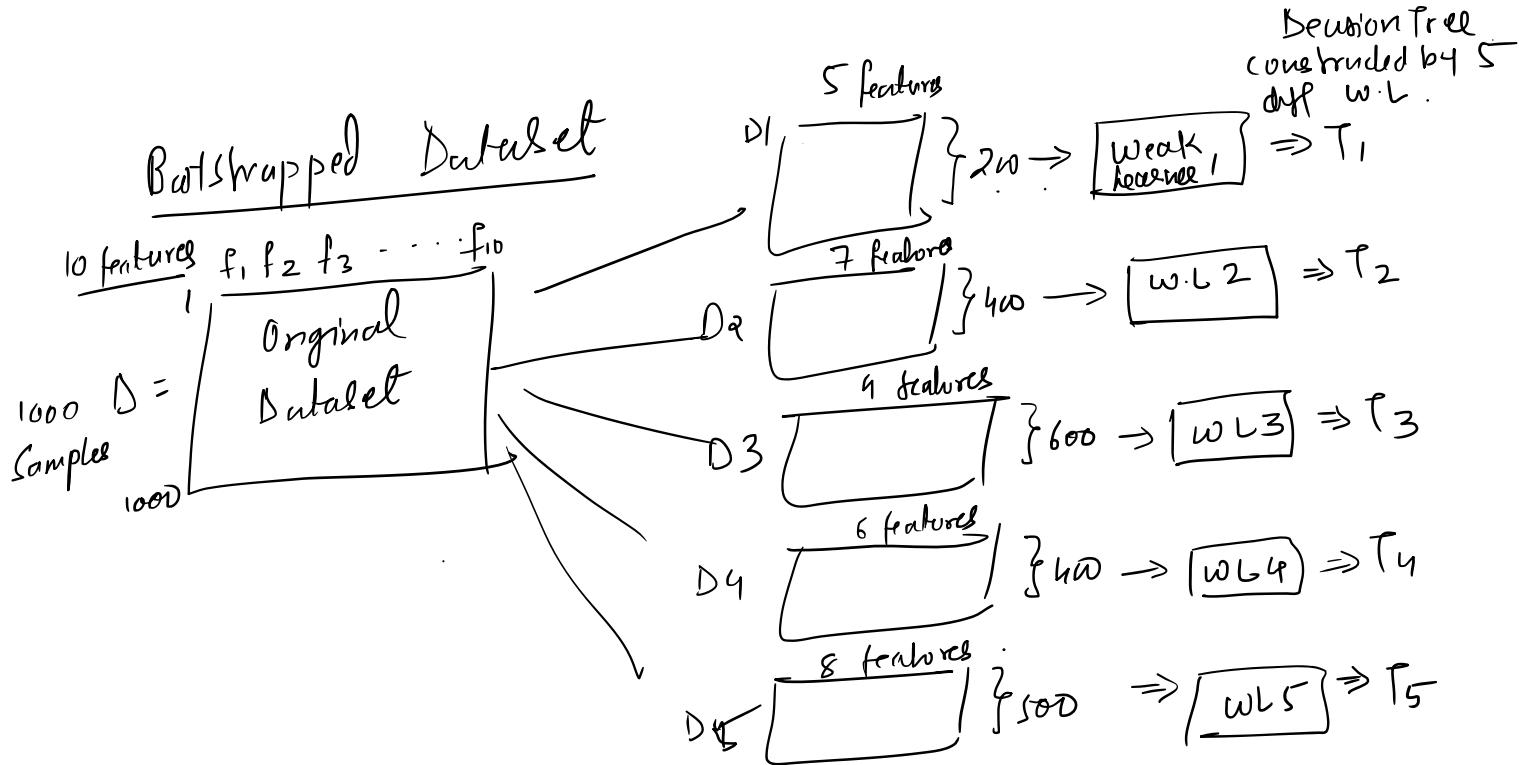
+ Decision Tree

- |
 - | Used for Regression \Rightarrow Regression Tree
 - | Used for Classification / Categorization \Rightarrow Classification Tree.

* Decision Tree is easy to Build and Interpret in practical.

* The Decision Tree is not very accurate on new test sample

* Random Forest Combines Simplicity of Decision Tree with flexibility resulting into vast improvement in accuracy.



We have 5 decision Tree Constructed one each by a weak learner

Now Test Sample = S_{test}

Now S_{test} is subjected to each of the 5 Decision Tree \Rightarrow

Suppose

$$T_1(S_{test}) = \text{Yes}$$

$$T_2(S_{test}) = \text{No}$$

$$T_3(S_{test}) = \text{Yes}$$

$$T_4(S_{test}) = \text{Yes}$$

$$T_5(S_{test}) = \text{No}$$

In 5 cases 3 cases are Yes
2 cases are No

Final Prediction Yes.

* Construct Random Forest

Original Set

	Chest Pain	Good Blood Circulation	Blocked Arteries	Weight	Heart Disease
S1	No	No	No	125	No
S2	Yes	Yes	Yes	180	Yes
S3	Yes	Yes	No	210	No
S4	Yes	No	Yes	167	Yes

No Yes Yes No ?

Step 1 → Create Bootstrapped Dataset

1. could be / not be of same size
2. Samples are randomly selected
3. Allowed to pick same sample more than once.

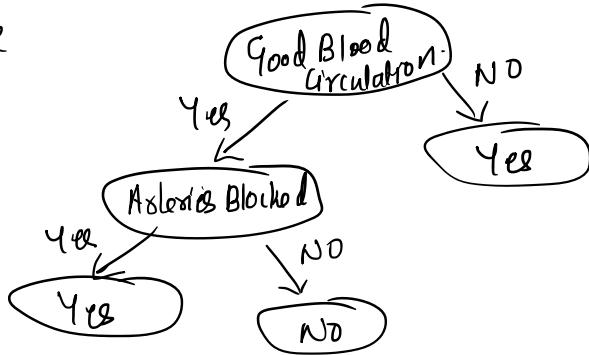
Bootstrapped Dataset

	Chest pain	Good Blood Circulation	Blocked Arteries	Weight	Heart Disease
	Yes	Yes	Yes	150	Yes
	No	No	No	125	No
	Yes	No	Yes	167	Yes
	Yes	No	Yes	167	Yes -

Step 2 → Create a Decision Tree Using Bootstrapped dataset but
Use random subset of Variables (features/columns)

Let us consider we use only 2 features
(Good Blood Circulation and Blocked Arteries)

Let the Tree be
(Gini Index)



Step 3 → go to step 1 and Repeat

- * Ideally we repeat it for 100 times
- * So we have T_1, T_2, \dots, T_{100} (large NO of Trees)
- * Each time the Tree Constructed is a weak learner.
[As all the features and samples are not considered while Tree construction].
- * Now we have Random Forest of 100 Trees and will be more effective than Individual Decision Tree.

⑥ How To Use the Random Forest (Here 100 Trees) →

Consider a Test Sample

Chest pain	Good Blood Circulation	Blocked Arteries	Weight	Heart Disease
Yes	No	No	168	?

Here we will run it through each of 100 Trees and will note the Prediction

Let say of 100 Trees

80 Trees Predicted Yes
20 Trees Predicted No

Answer = Yes

Answer = Yes

Q) To find Random Forest is Effective or Not ?

- * There might be some Sample not considered by any of Bootstrapped Dataset.
- * We will create a New Dataset with such samples. This is known as "Out of Bag" Dataset.
- * Now we will chk how many samples from Out of Bag Dataset are predicted correctly.
- * Numbers of Incorrectly Predicted Out of Bag Samples = Out of Bag Error.

Q) How to decide on Number of Columns to use for Building Trees in Random Forest ?

① In our Case DataSet has 4 features

→ Create a Random Forest of Trees using 2 features = F_1
→ Create a Random Forest of Trees using 3 features = F_2

100 Tree

100 Tree

② Now chk the Accuracy of Out of Bag Data on each of the Random Forest

③ Use the one that gives More Accurate Answer.

Summarize

Step 1: Create a Bootstrapped Data Set

Step 2: Create a Decision Tree using Subset of features

Step 3: Repeat Step 1 and 2 approx 100 times.

Step 4: Use Out of Bag Data Set to determine Out of Bag Error

Step 5: If Out of Bag Error is high then Repeat Step 1 to Step 4 until Out of Bag Error is considerably low.

Note → To compare the Performance for diff Random Forest Create diff Random Forest using different number of features

Use the one that gives Highest Accuracy.

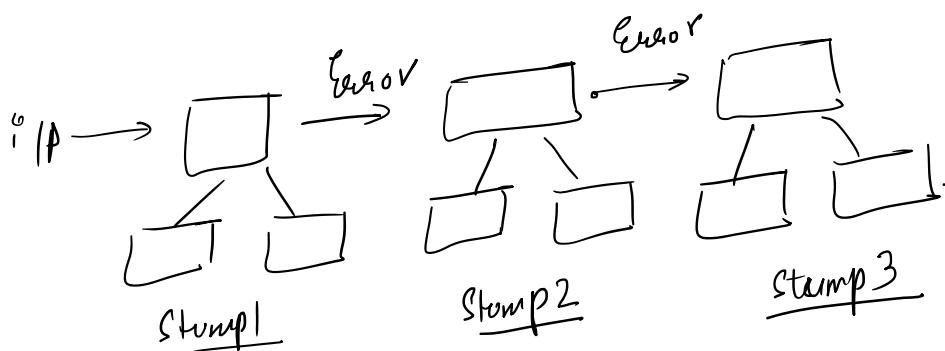
ADABOOST (ADAPTIVE BOOST)

→ (Tree Stumping)

- * AdaBoost creates Forest of Trees from scratch and then it uses it to make classification.
 - * [In Random Forest Everytime we make a Tree with different depth (depending on number of features considered).]
 - * There is no predetermined maximum depth of tree.]
 - * In ADABOOST, the forest of Trees are usually just node and Two leaves
 - * A Tree with one node and two leaves only is known as Tree Stump.
 - * In ADABOOST we have forest of Stump instead of Trees.
 - * Stump alone is not great for classification as it takes only one parameter to make decision.
 - * Stumps are weak learners.
 - * ADABOOST combines many stumps.
 - * In Random forest each tree contributes equally in final decision.

final decision

- * In ADABOOST, some stump may get more say in final decision than others.
- * In Random Forest the decision made by Trees are independent to other, so the order of tree generation is not important
- * In ADABOOST, we have forest of Stump and hence order is Important.
- * Here the error of first stump is corrected by next stump and so on -



→ Here Error of one stump influences the other stump
So the stumps are generated in sequence.

Idea Behind ADABOOST. (Boosting Technique \rightarrow reduce Bias)

1. Adaboost combines lot of weak learners to make classification
2. The weak learners are almost always Stump -
→ ... Some will have more say in classification

in the run

3. Some stump will have more say in classification than other stump.
4. Each stump is made by taking previous stump mistakes into account

⑥ Working of AdaBoost with Example.

Consider

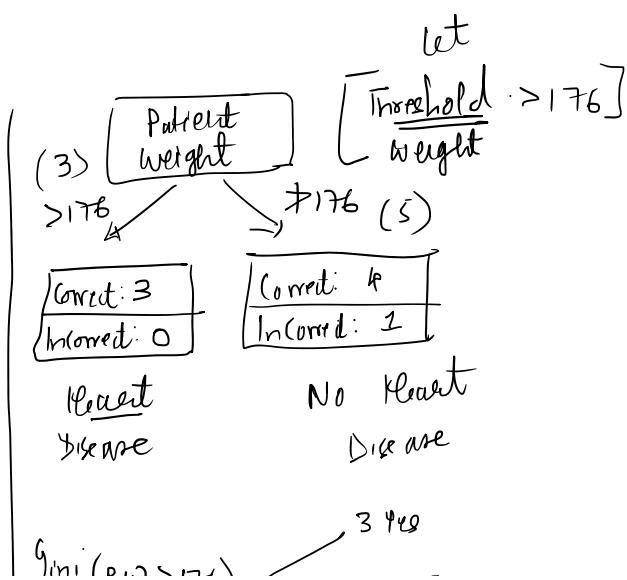
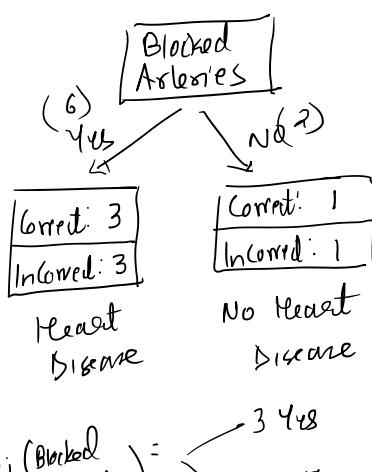
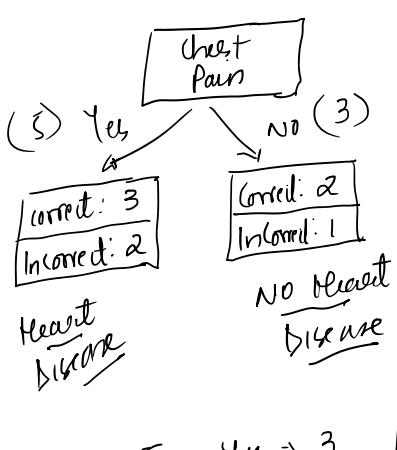
Chest Pain	Blocked Arteries	Patient Weight	Heart Disease
Yes	Yes	205	Yes
No	Yes	180	Yes
Yes	No	210	Yes
Yes	Yes	167	Yes
No	Yes	156	No
No	Yes	125	No
Yes	No	168	No
Yes	Yes	172	No

Step 1: Assign weight to every sample = $\frac{1}{\text{Total No of Sample}} = \frac{1}{8}$
(Initial)

* (Initially equal weights are assigned to each sample)

Chest Pain	Blocked Arteries	Patient Weight	Heart Disease	Sample weight
Yes	Yes	205	Yes	$\frac{1}{8}$
No	Yes	180	Yes	$\frac{1}{8}$
Yes	No	210	Yes	$\frac{1}{8}$
Yes	Yes	167	Yes	$\frac{1}{8}$
No	Yes	156	No	$\frac{1}{8}$
No	Yes	125	No	$\frac{1}{8}$
Yes	No	168	No	$\frac{1}{8}$
Yes	Yes	172	No	$\frac{1}{8}$

Step 2: Create stamp on each feature



$$Gini'(chest pain = \underline{Yes}) = \begin{cases} 4 \text{ Yes} \Rightarrow 3 \\ \text{No} \Rightarrow 2 \end{cases}$$

↑
Recent Disease

$$= 1 - ((\frac{3}{5})^2 + (\frac{2}{5})^2)$$

$$= \underline{0.48}$$

$$Gini'(chest pain = \underline{No}) = \begin{cases} 4 \text{ No} = 2 \\ 3 \text{ No} = 1 \end{cases}$$

$$= 1 - ((\frac{2}{3})^2 + (\frac{1}{3})^2)$$

$$= \underline{0.444}$$

weighted $Gini^0(\text{chest pain})$

$$= 0.48 * \frac{5}{8} + 0.44 * \frac{3}{8}$$

$$= \underline{\underline{0.4665}} \approx \underline{\underline{0.47}}$$

$Gini(\text{Blocked Arteries} = \underline{Yes}) = \begin{cases} 3 \text{ Yes} \\ 3 \text{ No} \end{cases}$ $= 1 - ((\frac{3}{6})^2 + (\frac{3}{6})^2)$ $= \underline{0.5}$	$Gini(\text{P.W} > 176) = \begin{cases} 3 \text{ Yes} \\ 0 \text{ No} \end{cases}$ $= 1 - ((\frac{3}{3})^2)$ $= \underline{0}$	$Gini(\text{P.W} \neq 176) = \begin{cases} 1 \text{ Yes} \\ 1 \text{ No} \end{cases}$ $= 1 - ((\frac{1}{5})^2 + (\frac{1}{5})^2)$ $= \underline{0.32}$
$Gini'(\text{Blocked Arteries}) = 0.5$	$Gini'(\text{weight}) = 0.2$	

* The smallest Gini Index is $\underline{0.2}$ for $\underline{\underline{\text{weight}}} > 176$ feature

So first Stump will be $\underline{\underline{\text{weight} > 176}}$

Step 3 → To calculate Amount of say of $\underline{\underline{\text{weight} > 176}}$ [feature].

$$\boxed{\text{Amount of Say} = \frac{1}{2} \ln \left(\frac{1 - \text{Total Error}}{\text{Total Error}} \right)}$$

Total Error: for a stamp the sum of weight associated with incorrectly classified sample.

Here only one sample is incorrectly classified.

Yes	Yes	167	Yes	18
-----	-----	-----	-----	----

$$\therefore \underline{\underline{\text{Total Error}}} = \underline{\underline{\frac{1}{8}}}$$

$$* \quad \text{Amt of Say} = \frac{1}{2} \ln \left(\frac{1 - \frac{1}{8}}{\frac{1}{8}} \right) = \frac{0.97}{(\text{High})}$$

[Note : If Amt of say is less then the stump is not doing good job]

- * now we have \rightarrow first stump (weight > 176)
 \rightarrow Amt of say.

Now initially every sample had equal weight \rightarrow

Note Before second stump we must modify weights of samples

- * Weights of correctly classified samples \Rightarrow must be decrease
- " " incorrectly classified samples \Rightarrow must be increase
- [So that next stump will focus more on incorrectly identified samples].

$$\begin{aligned} * \quad \text{New Sample weight (correctly classified)} &= \text{Sample weight (old)} \times e^{-\text{amount of say}} \\ \text{New Sample weight (incorrectly classified)} &= \text{Sample weight (old)} \times e^{+\text{amount of say}} \end{aligned}$$

$$\therefore \text{for Incorrectly Classified Sample} = \frac{1}{8} \times e^{+0.97} = 0.33$$

$$\text{For Correctly Classified Sample} = \frac{1}{8} \times e^{-0.97} = 0.05$$

Chest Pain	Blocked Arteries	Patient weight	Heart Disease	Sample weight	New Sample weight	Normalized Sample weight
Yes.	Yes	205	Yes	$\frac{1}{8}$	0.05	$0.05/0.68 = 0.07$
NO.	Yes	180	Yes	$\frac{1}{8}$	0.05	0.07
Yes.	No	210	Yes	$\frac{1}{8}$	0.05	0.07
Yes.	Yes	167	Yes	$\frac{1}{8}$	0.33	$0.33/0.68 = 0.49$

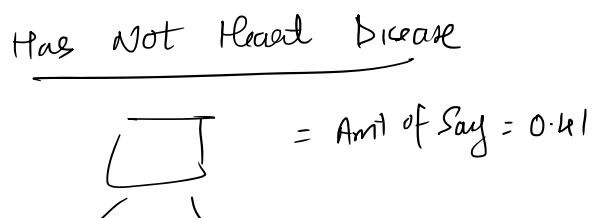
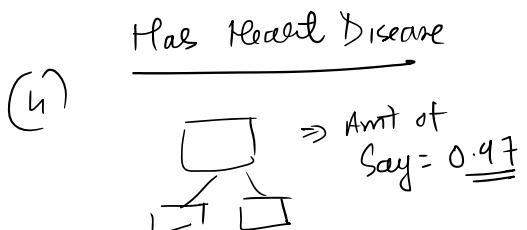
	No	210	Yes	$\frac{1}{8}$	0.05	
*	Yes	48	(167)	<u>Yes</u>	$\frac{1}{8} =$	<u>0.33</u>
✓	No	48	156	No	0.05	0.07
✓	No	48	125	No	0.05	0.07
.	Yes.	No	168	No	0.05	0.07
,	Yes	48	172	No	0.05	0.07
					$\uparrow \text{Sum} = 0.68$	

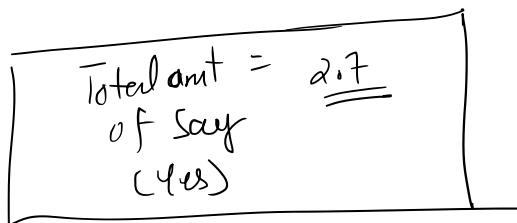
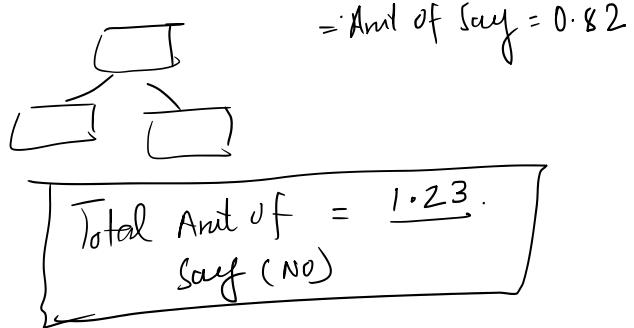
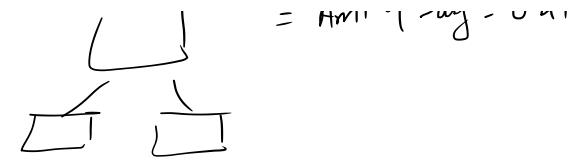
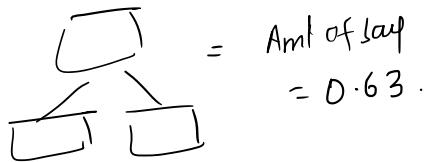
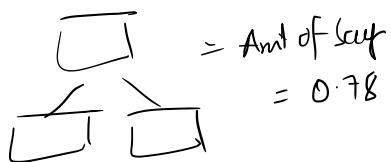
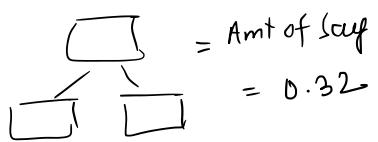
* Now we have New Normalized Sample weights:

- * The Incorrectly classified Sample has Higher New Sample weight
→ which means Stump 2 will focus more on it.
- * The Correctly classified Sample has Comparatively lower Sample weight
→ which means Stump 2 will focus less on it.
- * Similarly we can have Many such Stumps
- * This is how Adaboost creates and uses Stumps.

* Let us chk how forest of stumps made by Adaboost does Classification →

- ① let us consider a test sample
- ② we will pass the test sample through each stump
- ③ let there be few stumps that classifies Sample as Yes (^{Has Heart Disease}) and few stumps classifies Sample as No (^{Do not have heart disease}).





Total Amt of Say for stumps that classifies as Heart Disease
(Yes)

is greater

∴ The Test Sample is classified as Yes (Has Heart Disease).

Summary of Adaboost

- ① Assign sample weight (Initial equal)
- ② Create stump for each feature.
- ③ Use Gini Index to identify first stump.

- ④ For first stump:

- ① Calculate Total Error
- ② Amt of Say

④ Use "j" to train the stump.

⑤ For first stump:

① Calculate Total Error

② Amt of Say

⑥ Update the sample weights for each sample

$$\text{For Incorrectly classified weight} = \frac{\text{updated weight}}{\text{old weight}} + \text{Amt of Say}$$

$$\text{For correctly classified weight} = \frac{\text{updated weight}}{\text{old weight}} - \text{Amt of Say}$$

⑦ Now Specify updated weight for each sample & Normalize the weight to generate New Sample weight.

⑧ Identity new stump based on New Sample Weight & Repeat step 3 to 7

⑨ For Classification (Testing).

8.1 \Rightarrow Run the test sample through all the stumps.

8.2 \Rightarrow Calculate Amt of Say for sample classifying Yes & No

8.3 \Rightarrow Classify the test sample based on largest sum of Amount of Say