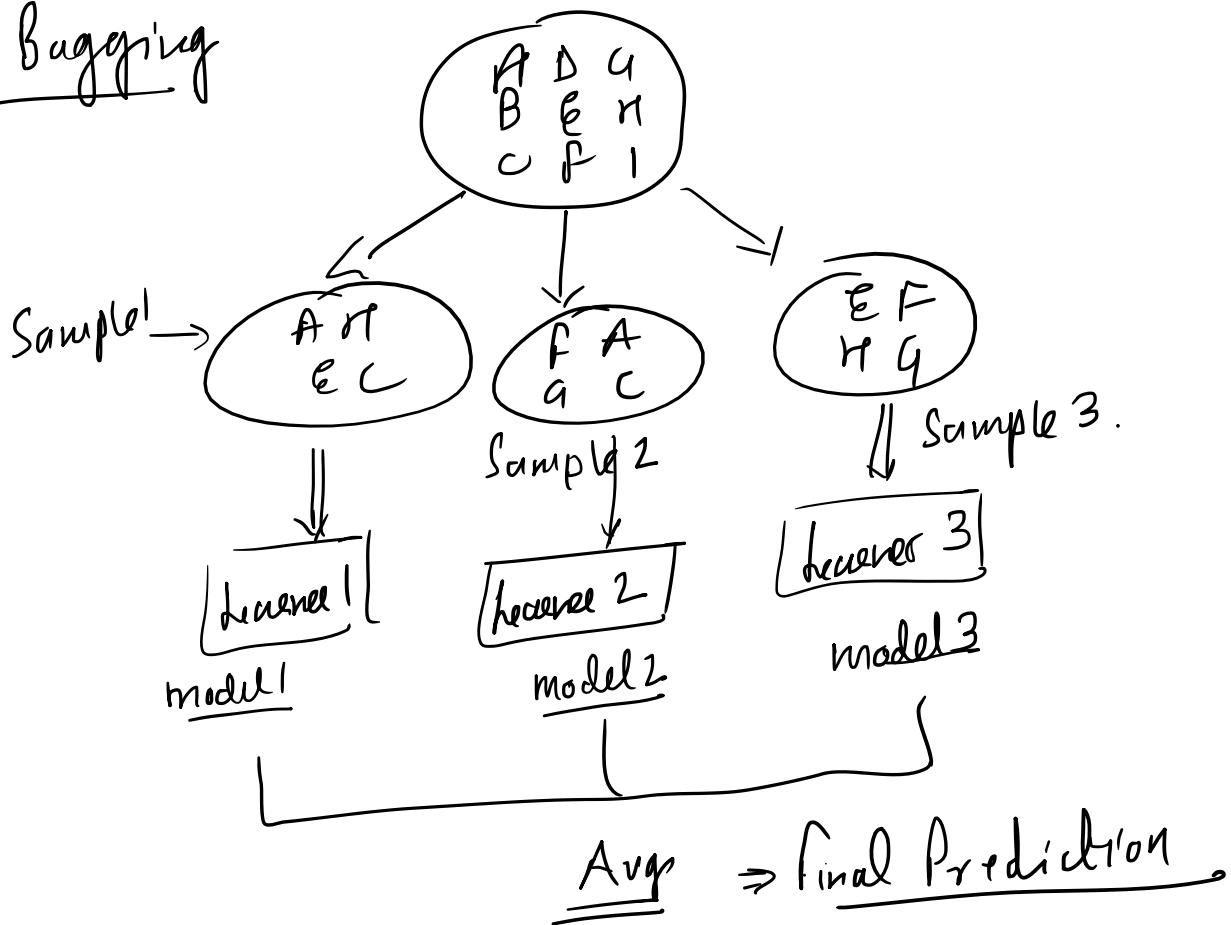


Note

Bagging

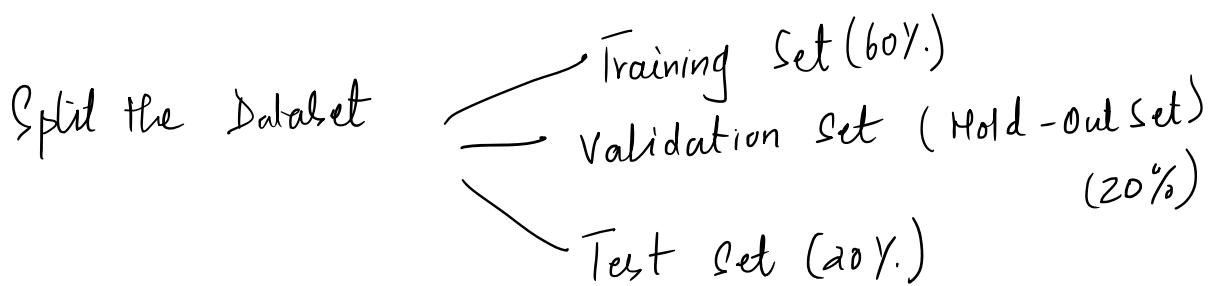


- * Bagging helps in reducing variance
- * Boosting helps in reducing bias.

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 ...
Validation set

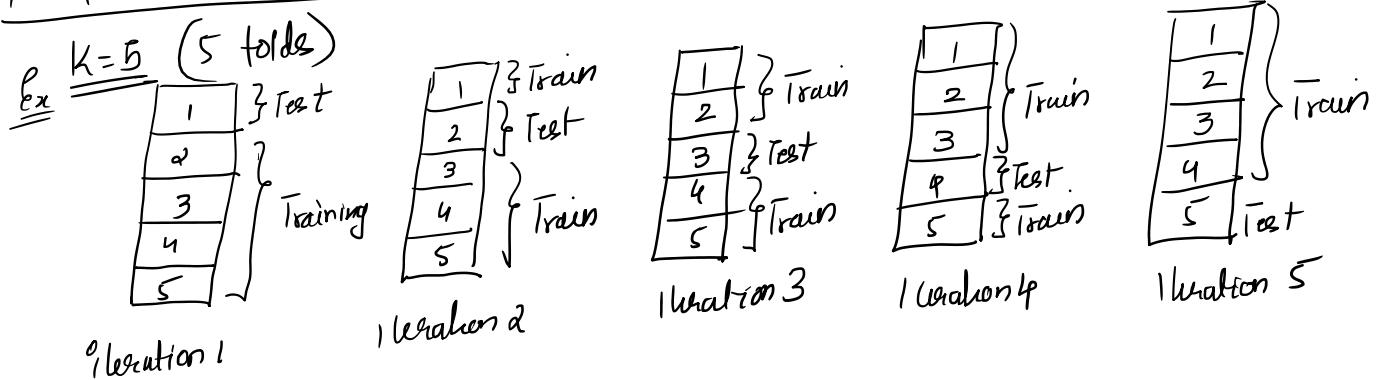
Test Set \Rightarrow^* Prediction

→ 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 →

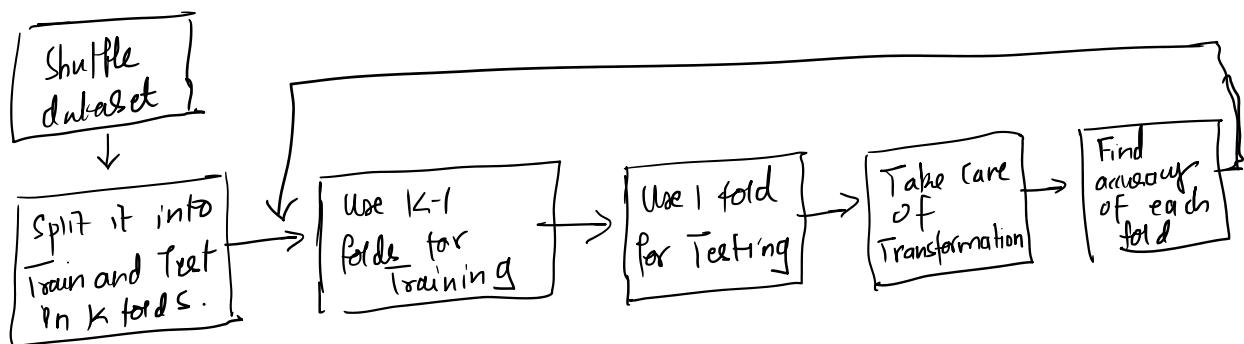


K fold → K fold helps us to build the model in generalized form.

→ 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

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

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 $\underline{k=10}$

\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]

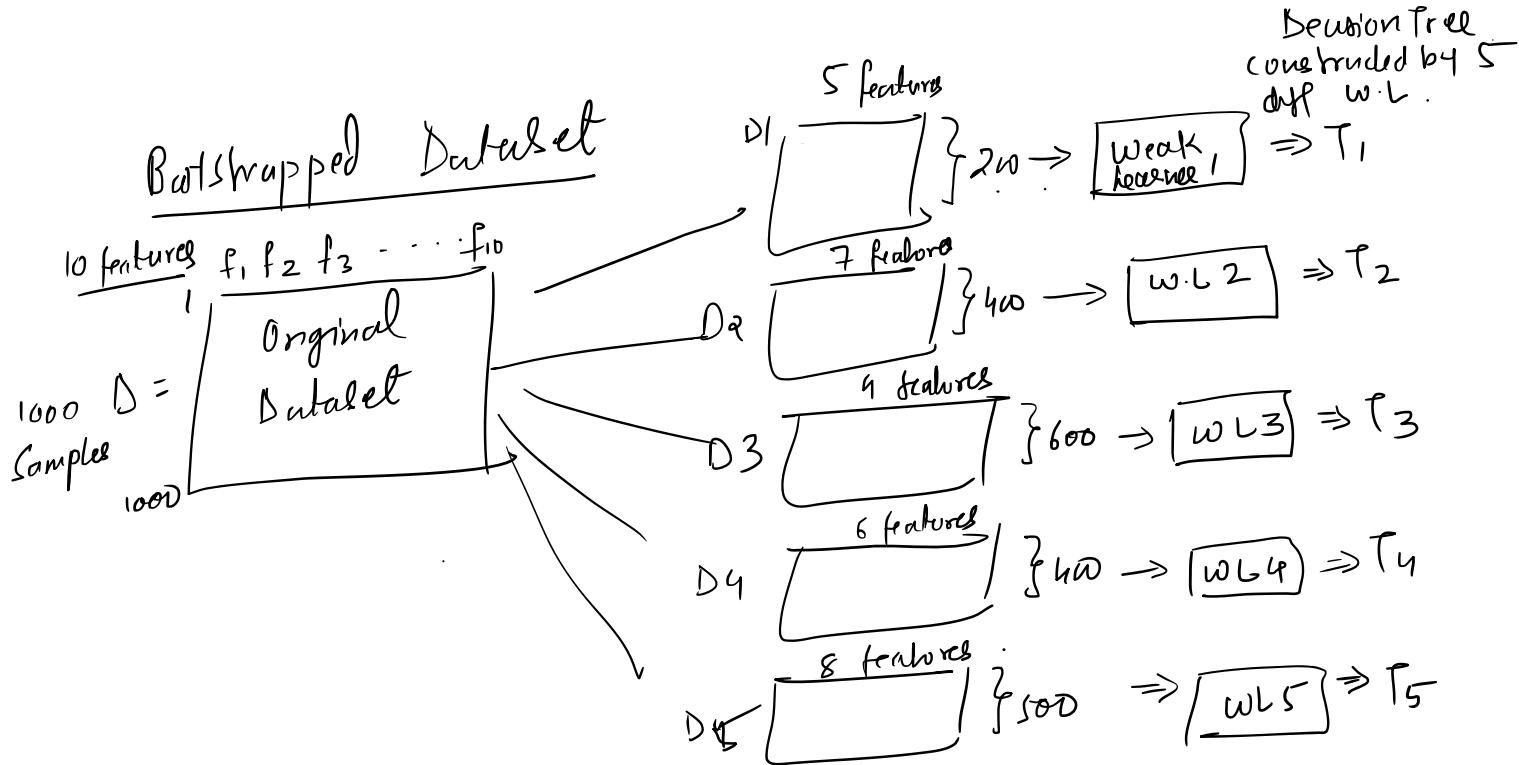
+ Decision Tree

- |
 - | 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 ⇒

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.

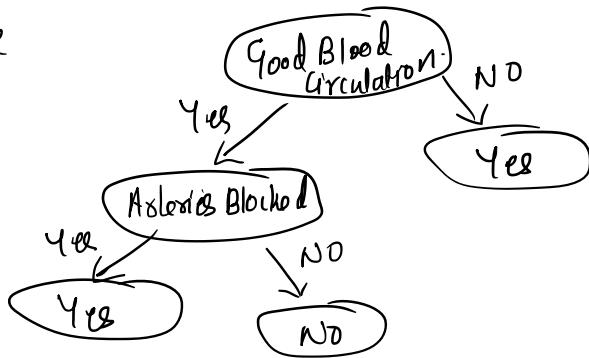
Bootstrapped Dataset

	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.

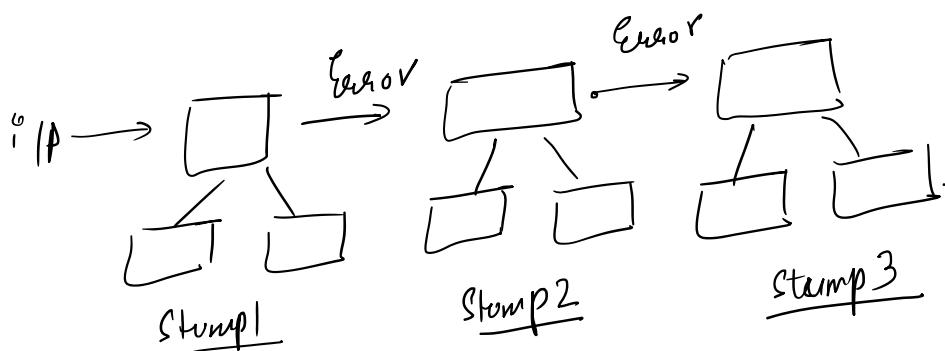
ADABOOST (ADAPTIVE BOOST)

→ (Tree Stomping)

- * 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